

The Zero-Error Feedback Capacity of State-Dependent Channels

Annina Bracher and Amos Lapidoth

July 11, 2016

Abstract

The zero-error feedback capacity of the Gelfand-Pinsker channel is established. It can be positive even if the channel's zero-error capacity is zero in the absence of feedback. Moreover, the error-free transmission of a single bit may require more than one channel use. These phenomena do not occur when the state is revealed to the transmitter causally, a case that is solved here using Shannon strategies. Cost constraints on the channel inputs or channel states are also discussed, as is the scenario where—in addition to the message—also the state sequence must be recovered.

1 Introduction

Motivated by Shannon's characterization of the zero-error capacity of the discrete memoryless channel (DMC) with a feedback link from the channel output to the encoder [1], we compute the corresponding capacity for the state-dependent DMC (SD-DMC) whose state is revealed acausally to the transmitter. This “Gelfand-Pinsker channel,” which was introduced by Gelfand and Pinsker in [2, 3], is more general than the channel studied by Shannon, and, indeed, when there is only one state we recover Shannon's result. But, more interestingly, this channel's zero-error feedback capacity exhibits phenomena that are not observed on the state-less channel: it can be positive even if the zero-error capacity is zero in the absence of feedback; the error-free transmission of a single bit may

The results in this paper were presented in part at the IEEE International Symposium on Information Theory (ISIT), Barcelona, Spain, Jul. 2016.

A. Bracher and A. Lapidoth are with the Department of Information Technology and Electrical Engineering, ETH Zurich, Switzerland (e-mail: bracher@isi.ee.ethz.ch; lapidoth@isi.ee.ethz.ch).

require more than one channel use; and Shannon’s sequential coding technique cannot be applied naively.

Like Shannon’s, our coding scheme is a two-phase scheme where the first phase reduces the receiver’s ambiguity to a manageable size, and the second removes it entirely. But our first phase differs from Shannon’s sequential approach and draws instead on Dueck’s scheme for zero-error communication over the multiple-access channel with feedback [4], which in turn draws on Ahlswede’s work [3, 5, 6]. The second phase is tricky, because sending a single bit reliably may require more than one channel use, so “uncoded” transmission need not work.

We also compute the zero-error feedback capacity of the SD-DMC $W(y|x, s)$ when the state is revealed to the transmitter causally. As we show, causal state information (SI) is utilized optimally using Shannon strategies. Consequently, when the SI is causal, the zero-error capacity is positive with feedback if, and only if, (iff) it is positive without it, and one channel use suffices to transmit a single bit error-free.

Several extensions are also discussed: we compute the zero-error feedback capacity of the Gelfand-Pinsker channel for the case where—in addition to the message—the encoder wishes to convey error-free also the state sequence; and we present capacity results for the Gelfand-Pinsker channel with cost constraints on the channel inputs or channel states. Under channel-input constraints a naive application of Shannon’s sequential coding technique turns out to be suboptimal even on the state-less channel.

The rest of this paper is structured as follows. We conclude this section by introducing some notation; by recalling the zero-error feedback capacity of the state-less DMC; and by exploring connections with the m -capacity of an arbitrarily-varying channel (AVC). Section 2 contains the problem formulation and the results. The main results for the Gelfand-Pinsker channel are proved in Section 3, and the paper concludes with a brief summary.

1.1 Notation and Terminology

We consider a SD-DMC of transition law $W(y|x, s)$, which is governed by an IID $\sim Q$ state process. The channel-input alphabet \mathcal{X} , the channel-state alphabet \mathcal{S} , and the channel-output alphabet \mathcal{Y} are all finite. By possibly redefining \mathcal{S} , we can assume without loss of generality that

$$Q(s) > 0, \quad s \in \mathcal{S}. \tag{1}$$

Subject to (1), the exact nature of the PMF Q is immaterial.

By default $\log(\cdot)$ denotes base-2 logarithm, and $\ln(\cdot)$ denotes natural logarithm. We denote by $h_b(\cdot)$ the binary entropy function. If ξ is a real number, then $[\xi]^+$ denotes the maximum of ξ and zero. Chance variables are denoted by upper-case letters and their realizations or the elements of their support sets by lower-case letters, e.g., Y denotes the random channel output and $y \in \mathcal{Y}$ a value it may take. Sets are denoted by calligraphic letters and in boldface if they are random, so the set of all messages is denoted \mathcal{M} , and \mathbf{M}_1 could be the set of messages of positive posterior probability given a first block of (random) channel outputs. Sequences are in bold lower- or upper-case letters depending on whether they are deterministic or random, e.g., \mathbf{Y} is the length- n channel-output sequence, and \mathbf{y} is an n -tuple from \mathcal{Y}^n . The positive integer $n \in \mathbb{N}$ stands for the blocklength, and unless otherwise specified sequences are of length n .

Variables pertaining to Time i have the subscript i , so S_i denotes the Time- i channel state. Sequences of variables that occur in the time-range j to i bear a subscript j and a superscript i , where the subscript $j = 1$ may be dropped, e.g., S_4^5 denotes the fourth and fifth state, and S^n denotes all the states through Time n . We also use a similar notation for sequences whose indices need not coincide with time, e.g., if \mathbf{s} is a 5-tuple from \mathcal{S}^5 , then s_3 denotes its third component, s_4^5 its fourth and fifth component, and s^5 the entire 5-tuple.

If the input X to the channel $W(y|x)$ is of PMF P , then $P \times W$ denotes the joint distribution of X and the channel output Y

$$(P \times W)(x, y) = P(x) W(y|x), \quad (x, y) \in \mathcal{X} \times \mathcal{Y},$$

and PW denotes the corresponding Y -marginal

$$(PW)(y) = \sum_{x \in \mathcal{X}} (P \times W)(x, y) = \sum_{x \in \mathcal{X}} P(x) W(y|x), \quad y \in \mathcal{Y}.$$

Given two PMFs P_1 and P_2 on some finite set \mathcal{Z} , we say that P_2 is absolutely continuous w.r.t. P_1 and write

$$P_2 \ll P_1,$$

if $P_2(z)$ is zero whenever $P_1(z)$ is. If P_2 is absolutely continuous w.r.t. P_1 , then the events that have probability zero w.r.t. P_1 must also have probability zero w.r.t. P_2 . Likewise for events of probability one.

For an SD-DMC $W(y|x, s)$ we denote by $\mathcal{P}(W)$ the set of transition laws $P_{Y|X, S}$ from $\mathcal{X} \times \mathcal{S}$ to \mathcal{Y} for which for every pair $(x, s) \in \mathcal{X} \times \mathcal{S}$

$$P_{Y|X, S}(\cdot|x, s) \ll W(\cdot|x, s).$$

For a state-less DMC $W(y|x)$ we drop s , and $\mathcal{P}(W)$ denotes the set of transition laws $P_{Y|X}$ from \mathcal{X} to \mathcal{Y} for which for every $x \in \mathcal{X}$

$$P_{Y|X}(\cdot|x) \ll W(\cdot|x).$$

The empirical type of an n -tuple $\mathbf{x} \in \mathcal{X}^n$ is denoted $P_{\mathbf{x}}$, i.e.,

$$P_{\mathbf{x}}(x) = \frac{N(x|\mathbf{x})}{n}, \quad x \in \mathcal{X},$$

where $N(x|\mathbf{x})$ is the number of components of the n -tuple \mathbf{x} that equal x . For a PMF P on \mathcal{X} the type class $\mathcal{T}_P^{(n)}$ comprises the elements of \mathcal{X}^n whose empirical type is P . If $\mathcal{T}_P^{(n)}$ is nonempty, then we say that P is an n -type. For an n -type P on \mathcal{X} , a transition law W from \mathcal{X} to \mathcal{Y} , and an element \mathbf{x} of $\mathcal{T}_P^{(n)}$ the W -shell $\mathcal{T}_W^{(n)}(\mathbf{x})$ comprises the n -tuples $\mathbf{y} \in \mathcal{T}_{P \times W}^{(n)}$ that satisfy $(\mathbf{x}, \mathbf{y}) \in \mathcal{T}_{P \times W}^{(n)}$.

1.2 State-Less Channels

Shannon showed in [1] that the zero-error capacity of the state-less DMC $W(y|x)$ (with or without feedback) is positive iff

$$\exists x, x' \in \mathcal{X} \text{ s.t. } \left(W(y|x) W(y|x') = 0, \forall y \in \mathcal{Y} \right). \quad (2)$$

When (2) holds, the error-free transmission of a single bit requires one channel use. He also showed that, when it is positive, the zero-error feedback capacity of $W(y|x)$ is

$$\max_{P_X} \min_{y \in \mathcal{Y}} -\log \sum_{x \in \mathcal{X}: W(y|x) > 0} P_X(x). \quad (3)$$

Ahlsweide [5] proved that (3) can be alternatively expressed as

$$\max_{P_X} \min_{P_{Y|X} \in \mathcal{P}(W)} I(X; Y), \quad (4)$$

where the mutual information is computed w.r.t. the joint PMF $P_X \times P_{Y|X}$. He also provided an alternative coding scheme. Unlike (2), the formulas (3) and (4) are only for channels with feedback. Indeed, feedback can increase the zero-error capacity of a DMC [1].

1.3 Connection to the AVC

There are interesting connections between the problem of computing the zero-error capacity of a DMC and that of computing the m -capacity (the capacity under the

maximal-probability-of-error criterion) of an AVC [7]. Indeed, given a DMC $W(y|x)$ with input alphabet \mathcal{X} and output alphabet \mathcal{Y} , the following construction produces an AVC $\widetilde{W}(y|x, \sigma)$ whose m-capacity is equal to the zero-error capacity of the channel $W(y|x)$ [7, Section 2], [8, Problem 12.3]. To construct the AVC we consider the functions $\sigma: \mathcal{X} \rightarrow \mathcal{Y}$ that satisfy that $W(\sigma(x)|x)$ is positive for all $x \in \mathcal{X}$. With each such function $\sigma(\cdot)$ we associate a state σ and the transition law

$$\widetilde{W}(y|x, \sigma) = \begin{cases} 1 & \text{if } y = \sigma(x), \\ 0 & \text{otherwise.} \end{cases} \quad (5)$$

The constructed AVC has two important properties. The first is that to every pair of input and output sequences x_1, \dots, x_n and y_1, \dots, y_n for which $\prod_k W(y_k|x_k)$ is positive, there corresponds a sequence of states $\sigma_1, \dots, \sigma_n$ such that $y_k = \sigma_k(x_k)$ for $k = 1, \dots, n$. The second is that $\widetilde{W}(y|x, \sigma)$ is $\{0, 1\}$ -valued in the sense that

$$\widetilde{W}(y|x, \sigma) \in \{0, 1\}, \quad \forall y, x, \sigma.$$

This latter property guarantees that the conditional probability of error over the AVC (conditional on the transmitted message and the state sequence) is $\{0, 1\}$ -valued and thus small (say, smaller than $1/2$) only if it is zero.

This relationship between the zero-error capacity and the m-capacity fails when the original channel whose zero-error capacity we seek is state-dependent and the state is revealed to the encoder. To see why, let us denote by $W(y|x, s)$ the transition law of the state-dependent channel whose zero-error capacity we seek when the state is revealed to the encoder, and suppose we want to construct an AVC $\widetilde{W}(y|x, \sigma)$ whose m-capacity when the state σ is revealed to the encoder is equal to the zero-error capacity we seek. We have intentionally used different letters s and σ for the state of the original channel and of the AVC because the two need not *prima facie* be the same. For example, if there is only one state s^* , then we are back to the state-less case and the construction we described above in (5) results in the number of AVC states being equal to the number of functions $\sigma: \mathcal{X} \rightarrow \mathcal{Y}$ that satisfy that $W(\sigma(x)|x, s^*)$ is positive for all $x \in \mathcal{X}$. However, in this case the m-capacity of the AVC $\widetilde{W}(y|x, \sigma)$ is equal to the zero-error capacity we seek only if the state σ is *not* revealed to the encoder. In attempting to construct the AVC we are faced with two conflicting requirements. For the state information (SI) that is revealed to the encoder in the two scenarios to be identical, the states s and σ should be identical. But for the AVC to have a $\{0, 1\}$ -law, the number of AVC states σ should typically be larger than the number of states s .

The construction does go through in the special case where the original state-dependent transition law $W(y|x, s)$ happens to be $\{0, 1\}$ -valued. In this special case

we can choose σ to equal s , and the m-capacity equals the zero-error capacity. In this case feedback is superfluous, because from the state (which is revealed to the encoder) and from the input (that it produces) the encoder can compute the output. We thus see that when $W(y|x, s)$ is $\{0, 1\}$ -valued the zero-error feedback capacity with acausal SI can be inferred from Ahlswede's results on the feedback-less AVC with SI at the encoder [9]; but in general it cannot.

2 Problem Formulation and Results

We consider an SD-DMC $W(y|x, s)$ with feedback whose encoder is furnished with the state sequence either acausally (Figure 1), or causally (Figure 2), or strictly-causally (Figure 4). Using n channel uses, the encoder wants to convey to the receiver error-free a message m from some finite set of messages \mathcal{M} . To this end it uses an (n, \mathcal{M}) zero-error code:

Definition 2.1. *Given a finite set \mathcal{M} and a positive integer $n \in \mathbb{N}$, an (n, \mathcal{M}) zero-error feedback code for the SD-DMC $W(y|x, s)$ with acausal SI to the encoder consists of n encoding mappings*

$$f_i: \mathcal{M} \times \mathcal{S}^n \times \mathcal{Y}^{i-1} \rightarrow \mathcal{X}, \quad i \in [1 : n] \quad (6)$$

and $|\mathcal{M}|$ disjoint decoding sets

$$\mathcal{D}_m \subseteq \mathcal{Y}^n, \quad m \in \mathcal{M}$$

such that, for every $m \in \mathcal{M}$ and every realization $\mathbf{s} \in \mathcal{S}^n$ of the state sequence, the probability of a decoding error is zero, i.e.,

$$\mathbb{P}[Y^n \notin \mathcal{D}_m | M = m, S^n = \mathbf{s}] = 0, \quad \forall m \in \mathcal{M}, \mathbf{s} \in \mathcal{S}^n,$$

where

$$\mathbb{P}[Y^n \notin \mathcal{D}_m | M = m, S^n = \mathbf{s}] = \sum_{\mathbf{y} \in \mathcal{Y}^n \setminus \mathcal{D}_m} \prod_{i=1}^n W(y_i | f_i(m, \mathbf{s}, y^{i-1}), s_i). \quad (7)$$

A rate R is achievable if for every sufficiently-large blocklength n there exists an (n, \mathcal{M}) zero-error feedback code with

$$\log |\mathcal{M}| \geq nR.$$

The zero-error feedback capacity with acausal SI is the supremum of all achievable rates and is denoted $C_{f,0}$.

The zero-error feedback capacities with causal and strictly-causal SI are denoted $C_{f,0}^{\text{caus}}$ and $C_{f,0}^{\text{s-caus}}$, respectively. They are defined like $C_{f,0}$ except that the encoding mappings (6) are replaced by

$$f_i: \mathcal{M} \times \mathcal{S}^i \times \mathcal{Y}^{i-1} \rightarrow \mathcal{X}, \quad i \in [1 : n] \quad (8)$$

in the causal case and by

$$f_i: \mathcal{M} \times \mathcal{S}^{i-1} \times \mathcal{Y}^{i-1} \rightarrow \mathcal{X}, \quad i \in [1 : n] \quad (9)$$

in the strictly-causal case.

Note that the PMF Q governing the state does not appear in Definition 2.1 and therefore does not affect the zero-error feedback capacities with acausal, causal, and strictly-causal SI. Also note that our definition assumes deterministic encoders. This assumption is not restrictive:

Remark 2.2. *Allowing stochastic encoders does not increase the zero-error feedback capacities with acausal, causal, and strictly-causal SI.*

Proof. A proof for the case where the encoder observes the SI acausally is provided in Appendix A. The proof goes through also when the SI is causal or strictly-causal. \square

2.1 Acausal SI

In this section we assume that the encoder observes the SI acausally (see Figure 1). Our main result is presented in the following two theorems, which together provide a single-letter characterization of $C_{f,0}$. The first characterizes the channels for which it is positive, and the second provides a formula for $C_{f,0}$ when it is positive.

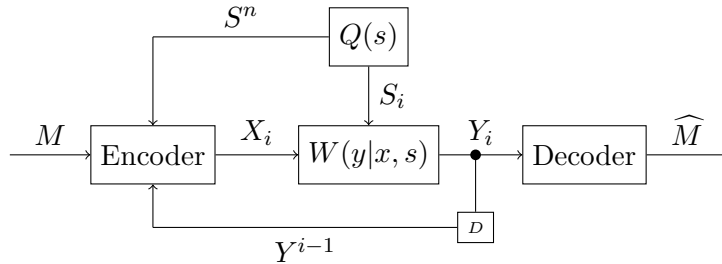


Figure 1: SD-DMC with acausal SI and feedback.

Theorem 2.3. *A necessary and sufficient condition for $C_{f,0}$ to be positive is*

$$\forall s, s' \in \mathcal{S} \quad \exists x, x' \in \mathcal{X} \text{ s.t. } (W(y|x, s) W(y|x', s') = 0, \forall y \in \mathcal{Y}). \quad (10)$$

Proof. See Section 3.1. \square

Theorem 2.4. *If $C_{f,0}$ is positive, then*

$$C_{f,0} = \min_{P_S} \max_{P_{U,X|S}} \min_{\substack{P_{Y|U,X,S}: \\ P_{Y|U=u,X,S} \in \mathcal{P}(W), \forall u \in \mathcal{U}}} I(U; Y) - I(U; S), \quad (11)$$

where U is an auxiliary chance variable taking values in a finite set \mathcal{U} , and the mutual informations are computed w.r.t. the joint PMF $P_S \times P_{U,X|S} \times P_{Y|U,X,S}$. Restricting X to be a function of U and S , i.e., $P_{U,X|S}$ to have the form

$$P_{U,X|S}(u, x|s) = P_{U|S}(u|s) \mathbb{1}_{x=g(u,s)}, \quad (12)$$

does not change the RHS of (11), nor does restricting the cardinality of \mathcal{U} to

$$|\mathcal{U}| \leq |\mathcal{X}|^{|\mathcal{S}|}. \quad (13)$$

Proof. See Section 3.2. \square

Remark 2.5. *The hypothesis in Theorem 2.4 that $C_{f,0}$ be positive is essential: the RHS of (11) may be positive even when $C_{f,0}$ is zero.*

In fact, as we prove in Appendix B:

Remark 2.6. *The RHS of (11) is positive iff*

$$\forall (s, y) \in \mathcal{S} \times \mathcal{Y} \quad \exists x \in \mathcal{X} \text{ s.t. } W(y|x, s) = 0. \quad (14)$$

Theorems 2.3 and 2.4 generalize to the SD-DMC with feedback and acausal SI Shannon's characterization [1, Theorem 7] of the zero-error feedback capacity of the (state-less) DMC $W(y|x)$ (see (2) and (3) in Section 1.2). That (10) reduces to (2) when $|\mathcal{S}| = 1$ is evident. That (11) reduces to (3) when $|\mathcal{S}| = 1$ becomes evident when we recall from [5] Ahlswede's alternative form (4) for (3): clearly, (11) specializes to (4) and thus to (3) when $|\mathcal{S}| = 1$. The way in which (11) generalizes (4) is reminiscent of the way the Gelfand-Pinsker capacity generalizes the ordinary capacity of the state-less DMC (cf. [2, 10]).

In the remainder of this section we discuss how feedback affects the zero-error capacity with acausal SI. By considering the case of a single state, i.e., $|\mathcal{S}| = 1$, and invoking

Shannon’s result [1] that feedback can increase the zero-error capacity of a DMC, we readily obtain that feedback can also increase the zero-error capacity of an SD-DMC with acausal SI. But, in the presence of acausal SI, more is true. Unlike the stateless channel, here feedback can increase the capacity from zero:

Theorem 2.7. *The zero-error capacity of an SD-DMC with acausal SI can be positive with feedback yet zero without it.*

Proof. See Section 3.3. □

Condition (10) is thus only for channels with feedback: the no-feedback zero-error capacity of the SD-DMC $W(y|x, s)$ with acausal SI can be zero also when the channel satisfies (10). Because feedback can help only if the encoder uses the channel more than once, we obtain the following corollary, which marks another difference to the state-less case:

Corollary 2.8. *On the SD-DMC with acausal SI and feedback, the error-free transmission of a single bit may require more than one channel use.*

This result will be strengthened in Section 2.2, where we show that also in the absence of feedback the error-free transmission of a single bit may require more than one channel use (Corollary 2.15).

As we have seen in Section 1.3, if the transition law $W(y|x, s)$ of the SD-DMC happens to be $\{0, 1\}$ -valued, then $C_{f,0}$ is related to Ahlswede’s AVC with acausal SI. As we show in Appendix C, in this case Theorems 2.3 and 2.4 can be greatly simplified:

Example 2.9. *If the transition law $W(y|x, s)$ of an SD-DMC is $\{0, 1\}$ -valued, then*

$$C_{f,0} = \min_{s \in \mathcal{S}} \log |\{y \in \mathcal{Y} : \exists x \in \mathcal{X} \text{ s.t. } W(y|x, s) > 0\}|. \quad (15)$$

Remark 2.5 notwithstanding, if $W(y|x, s)$ is $\{0, 1\}$ -valued, then the RHS of (11)—which in this case is equal to the RHS of (15)—is positive iff $C_{f,0}$ is positive. This agrees with Ahlswede’s observation [9] that the formula for the (a- and m-) capacity of the general AVC $W(y|x, s)$ whose state sequence is revealed acausally to the encoder not only applies when the capacity is positive but also determines whether it is positive.

2.2 Causal SI

In this section we assume that the encoder observes the SI causally (see Figure 2). The following two theorems together provide a single-letter characterization of $C_{f,0}^{\text{caus}}$. The

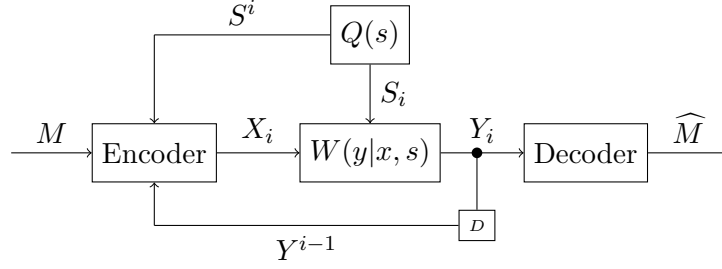


Figure 2: SD-DMC with causal SI and feedback.

first characterizes the channels for which it is positive, and the second provides a formula for the capacity when it is positive.

Theorem 2.10. *A necessary and sufficient condition for $C_{f,0}^{\text{caus}}$ to be positive is that there exist a partition $\mathcal{Y}_0, \mathcal{Y}_1$ of \mathcal{Y} for which*

$$\forall s \in \mathcal{S} \quad \exists x, x' \in \mathcal{X} \text{ s.t. } W(\mathcal{Y}_0|x, s) = W(\mathcal{Y}_1|x', s) = 1. \quad (16)$$

If $C_{f,0}^{\text{caus}}$ is positive, then one channel use suffices to transmit a single bit error-free, and therefore the zero-error capacity with causal SI is positive with feedback iff it is positive without it.

Proof. See Appendix F. □

Theorem 2.11. *If $C_{f,0}^{\text{caus}}$ is positive, then*

$$C_{f,0}^{\text{caus}} = \max_{P_U} \min_{P_{Y|U} \in \mathcal{P}(W')} I(U; Y) \quad (17)$$

$$= \max_{P_U} \min_y -\log \sum_{u: W'(y|u) > 0} P_U(u), \quad (18)$$

where U is an auxiliary chance variable taking values in a finite set \mathcal{U} of cardinality $|\mathcal{U}| = |\mathcal{X}|^{|\mathcal{S}|}$; the mutual information is computed w.r.t. the joint PMF $P_U \times P_{Y|U}$; and

$$W'(y|u) = \sum_{s \in \mathcal{S}} Q_S(s) W(y|g(u, s), s), \quad (u, y) \in \mathcal{U} \times \mathcal{Y}, \quad (19)$$

where $\{g(u, \cdot) : u \in \mathcal{U}\}$ is the set of functions from \mathcal{S} to \mathcal{X} , i.e., $\mathcal{X}^{\mathcal{S}}$. Because

$$(W'(y|u) > 0) \iff (\exists s \in \mathcal{S} \text{ s.t. } W(y|g(u, s), s) > 0), \quad (20)$$

$P_{Y|U} \in \mathcal{P}(W')$ holds iff

$$(W(y|g(u, s), s) = 0, \forall s \in \mathcal{S}) \implies (P_{Y|U}(y|u) = 0). \quad (21)$$

Proof. The proof draws on Shannon's results [1, 11] (see Appendix G). \square

Remark 2.12. *The hypothesis in Theorem 2.4 that $C_{f,0}^{\text{caus}}$ be positive is essential: the RHS of (17) may be positive even when $C_{f,0}^{\text{caus}}$ is zero.*

In fact, as we prove in Appendix H:¹

Remark 2.13. *The RHS of (17) is positive iff*

$$\forall (s, y) \in \mathcal{S} \times \mathcal{Y} \quad \exists x \in \mathcal{X} \text{ s.t. } W(y|x, s) = 0. \quad (22)$$

Theorems 2.10 and 2.11 generalize to the SD-DMC with feedback and causal SI Shannon's characterization [1, Theorem 7] of the zero-error feedback capacity of the (state-less) DMC $W(y|x)$ (see (2) and (3) in Section 1.2). The way in which (16) and (18) generalize (2) and (3) is reminiscent of the way the ordinary capacity with causal SI generalizes the ordinary capacity of the state-less DMC (cf. [10, 11]): in both cases causal SI is utilized optimally by using *Shannon strategies*. To see this, recall that by using Shannon strategies the encoder transforms the SD-DMC $W(y|x, s)$ with causal SI into the state-less DMC

$$W'(y|u) = \sum_{s \in \mathcal{S}} Q_S(s) W(y|g(u, s), s)$$

with input alphabet \mathcal{U} of cardinality $|\mathcal{U}| = |\mathcal{X}|^{|\mathcal{S}|}$, where $\{g(u, \cdot) : u \in \mathcal{U}\}$ equals $\mathcal{X}^{\mathcal{S}}$: an encoder with causal SI is said to use Shannon strategies if it performs the encoding over the set \mathcal{U} and obtains the Time- i channel-input by evaluating the function $g(\cdot, \cdot) : \mathcal{U} \times \mathcal{S} \rightarrow \mathcal{X}$ for the i -th codeword-symbol $u_i \in \mathcal{U}$ and the Time- i channel-state S_i (see Figure 3 and [12, Remark 7.6]). By comparing (16) and (18) to (2) and (3), respectively, we see that, indeed, the zero-error feedback capacity of the SD-DMC $W(y|x, s)$ with causal SI equals the zero-error feedback capacity of the state-less DMC $W'(y|u)$, and hence causal SI is utilized optimally by using Shannon strategies.

In the remainder of this section we briefly contrast how feedback affects the zero-error capacities with acausal and causal SI. As in the acausal case, by considering the case of a single state, i.e., $|\mathcal{S}| = 1$, and invoking Shannon's result [1] that feedback can increase the zero-error capacity of a DMC, we readily obtain that feedback can also

¹Remarks 2.6 and 2.13 imply that, like the ordinary capacities with causal and acausal SI [2, 11], the RHS of (11) is positive iff that of (17) is positive. As we shall see, however, this does not hold for the capacities: the zero-error capacity can be positive with acausal SI yet zero with causal SI (see Theorem 2.14 ahead).

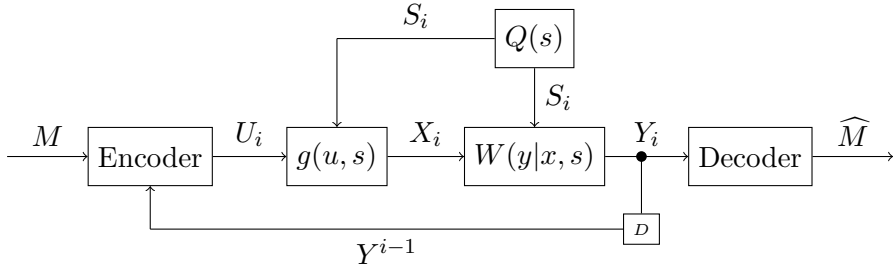


Figure 3: Shannon strategies.

increase the zero-error capacity in the causal case. However, unlike the acausal case, the zero-error capacity with causal SI is positive with feedback iff it is positive without it (see Theorem 2.10).

Since acausal SI is better than causal SI, and since the zero-error capacity with causal SI is positive with feedback iff it is positive without it, the condition in Theorem 2.10 is sufficient for the no-feedback zero-error capacity of the SD-DMC $W(y|x, s)$ with acausal SI to be positive. (Alternatively, this is obtained by noting that (16) of Theorem 2.10 implies (143) of Lemma 3.3 ahead, which is a sufficient condition for the no-feedback zero-error capacity of the SD-DMC with acausal SI to be positive. As we shall see in Example 2.16 ahead, the reverse implication need not hold.) By Theorem 2.7 the zero-error capacity with acausal SI can be positive with feedback yet zero without it. Consequently, unlike the ordinary capacities with causal and acausal SI (see [2, 11]) or the RHSs of (11) and (17) (see Remarks 2.6 and 2.13), the zero-error feedback capacity can be positive with acausal SI yet zero with causal SI. In fact, more is true:

Theorem 2.14. *The zero-error capacity can be positive with acausal SI yet zero with causal SI even when feedback is available in the latter setting and absent in the former.*

Because acausal SI can be better than causal SI only if the encoder uses the channel more than once, we obtain the following corollary, which strengthens Corollary 2.8:

Corollary 2.15. *On the SD-DMC with acausal SI, the error-free transmission of a single bit may require more than one channel use also in the absence of feedback.*

To prove Theorem 2.14, we provide an example for which the zero-error capacity (with and without feedback) is positive with acausal SI yet zero with causal SI:

Example 2.16. *Consider a deterministic SD-DMC $W(y|x, s)$ over the alphabets $\mathcal{X} = \{0, 1\}$ and $\mathcal{S} = \mathcal{Y} = \{1, 2, 3\}$. Let the output corresponding to the input x and the state*

s be the single element of the set $\mathcal{Y}_{x,s}$ that is given in Table 1

$$\{y \in \mathcal{Y} : W(y|x, s) > 0\} = \mathcal{Y}_{x,s}, \quad \forall (x, s) \in \mathcal{X} \times \mathcal{S}. \quad (23)$$

Since this channel violates (16) but satisfies (143) of Lemma 3.3 ahead for $\kappa = \lambda = 3$,

$$x(s, k) = \begin{cases} 0 & \text{if } k = 1 \text{ or } (s, k) = (3, 2), \\ 1 & \text{otherwise,} \end{cases} \quad (s, k) \in \{1, 2, 3\} \times \{1, 2, 3\},$$

and $\mathcal{Y}_\ell = \{\ell\}$, $\ell \in \{1, 2, 3\}$ (cf. Remark 3.5 ahead), its zero-error capacity (both with and without feedback) is positive with acausal SI yet zero with causal SI.

$\mathcal{Y}_{x,s}$		s		
		1	2	3
x	0	$\{2\}$	$\{1\}$	$\{1\}$
	1	$\{3\}$	$\{3\}$	$\{2\}$

Table 1: Nonzero transitions of the SD-DMC in Example 2.16.

2.3 Strictly-Causal SI

In this section we assume that the encoder observes the SI strictly-causally (see Figure 4).

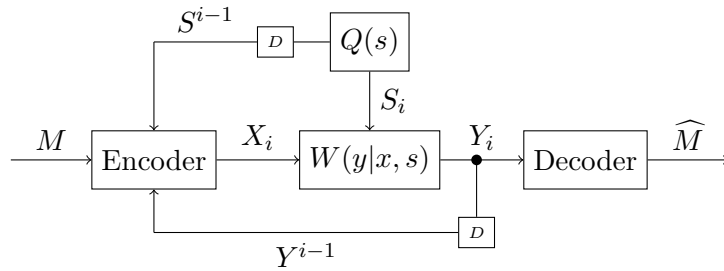


Figure 4: SD-DMC with strictly-causal SI and feedback.

The results (2)–(4) for the state-less DMC also provide the zero-error feedback capacity $C_{f,0}^{\text{s-causal}}$ of the state-dependent channel whose state is revealed strictly-causally to the encoder:

Remark 2.17. Shannon’s proof of (2) and (3) in [1] goes through also when the channel is state-dependent and the SI is revealed strictly-causally to the encoder. Consequently, such SI cannot increase the zero-error feedback capacity. That is, if we define

$$\widetilde{W}(y|x) = \sum_{s \in \mathcal{S}} Q(s) W(y|x, s), \quad (x, y) \in \mathcal{X} \times \mathcal{Y}, \quad (24)$$

then a necessary and sufficient condition for $C_{f,0}^{\text{sc-caus}}$ to be positive is that (2) hold for the channel $\widetilde{W}(y|x)$, and if $C_{f,0}^{\text{sc-caus}}$ is positive, then it can be computed by substituting $\widetilde{W}(y|x)$ for $W(y|x)$ in (3) or (4).²

2.4 Zero-Error Rate-and-State

In this section we consider a scenario where—in addition to the message m —the encoder wishes to convey to the receiver (error-free) also the state sequence S^n , which it observes acausally. For the standard setting where the probability of a message error need not be zero but can be arbitrarily small, Kim, Sutivong, and Cover [13] introduced and solved a related problem with list decoding of state sequences. Choudhuri, Kim, and Mitra [14] studied the causal and strictly-causal settings subject to a constraint on the distortion between the state sequence and its receiver-side estimate. Analogous results in the presence of feedback were recently reported by Bross and Lapidoth [15].

We begin with the basic definitions of an (n, \mathcal{M}) zero-error code:

Definition 2.18. Given a finite set \mathcal{M} and a positive integer $n \in \mathbb{N}$, an (n, \mathcal{M}) zero-error state-conveying feedback code for the SD-DMC $W(y|x, s)$ with acausal SI to the encoder consists of n encoding mappings

$$f_i: \mathcal{M} \times \mathcal{S}^n \times \mathcal{Y}^{i-1} \rightarrow \mathcal{X}, \quad i \in [1 : n]$$

and $|\mathcal{M}| |\mathcal{S}|^n$ disjoint decoding sets

$$\mathcal{D}_{m,\mathbf{s}} \subseteq \mathcal{Y}^n, \quad (m, \mathbf{s}) \in \mathcal{M} \times \mathcal{S}^n$$

such that for every $(m, \mathbf{s}) \in \mathcal{M} \times \mathcal{S}^n$ the probability of a decoding error is zero, i.e.,

$$\mathbb{P}[Y^n \notin \mathcal{D}_{m,\mathbf{s}} | M = m, S^n = \mathbf{s}] = 0, \quad \forall (m, \mathbf{s}) \in \mathcal{M} \times \mathcal{S}^n,$$

where

$$\mathbb{P}[Y^n \notin \mathcal{D}_{m,\mathbf{s}} | M = m, S^n = \mathbf{s}] = \sum_{\mathbf{y} \in \mathcal{Y}^n \setminus \mathcal{D}_{m,\mathbf{s}}} \prod_{i=1}^n W(y_i | f_i(m, \mathbf{s}, y^{i-1}), s_i). \quad (25)$$

²Note that by (1) $\widetilde{W}(y|x)$ is positive iff there exists some state for which $W(y|x, s)$ is positive.

A positive rate R is called *achievable* if for every sufficiently-large blocklength n there exists an (n, \mathcal{M}) zero-error state-conveying feedback code satisfying

$$\frac{1}{n} \log |\mathcal{M}| \geq R.$$

The zero-error state-conveying feedback capacity is the supremum of all achievable rates and is denoted $C_{f,0}^{m+s}$. If no positive rate is achievable, then we say that $C_{f,0}^{m+s} = 0$, regardless of whether or not it is possible to convey the state sequence error-free.

Our definition of an (n, \mathcal{M}) zero-error state-conveying code does not depend on the PMF Q of the state and assumes a deterministic encoder. Like the scenario where the encoder need not convey the state, $C_{f,0}^{m+s}$ does not depend on the PMF Q , and allowing stochastic encoders does not increase it (cf. the proof of Remark 2.2).

The following theorem provides a single-letter characterization of $C_{f,0}^{m+s}$:

Theorem 2.19. *A necessary condition for $C_{f,0}^{m+s}$ to be positive is (10), and if (10) holds, then*

$$C_{f,0}^{m+s} = \left[\min_{P_S} \max_{P_{X|S}} \min_{P_{Y|X,S} \in \mathcal{P}(W)} I(X, S; Y) - H(S) \right]^+, \quad (26)$$

where the mutual information and the entropy are computed w.r.t. the joint PMF $P_S \times P_{X|S} \times P_{Y|X,S}$.

Proof. The result is proved in Appendix I by adapting the proofs of Theorems 2.3 and 2.4 so as to guarantee that the receiver can decode also the state sequence error-free. \square

2.5 Constrained Inputs

In this section we establish the zero-error feedback capacity of the SD-DMC $W(y|x, s)$ with acausal SI subject to a cost constraint on the channel inputs. Consider some nonnegative “cost-function” $\gamma: \mathcal{X} \rightarrow \mathbb{R}_0^+$, and define

$$\gamma_{\min} = \min_{x \in \mathcal{X}} \gamma(x) \quad \text{and} \quad \gamma_{\max} = \max_{x \in \mathcal{X}} \gamma(x).$$

Let the set $\mathcal{X}' \subseteq \mathcal{X}$ comprise all the minimizers of $\gamma(\cdot)$

$$\mathcal{X}' = \{x \in \mathcal{X}: \gamma(x) = \gamma_{\min}\}.$$

The cost constraint we study is that, at every blocklength n and for every transmitted message $m \in \mathcal{M}$, the channel inputs’ average cost

$$\gamma^{(n)}(X^n) = \frac{1}{n} \sum_{i=1}^n \gamma(X_i)$$

satisfy the cost constraint

$$\gamma^{(n)}(X^n) \leq \Gamma \quad (27)$$

for some given Γ satisfying

$$\gamma_{\min} < \Gamma < \gamma_{\max}. \quad (28)$$

The zero-error feedback capacity with acausal SI subject to (27) is denoted $C_{f,0}(\Gamma)$. We restrict Γ to (28), because all other values of Γ are uninteresting: if $\Gamma < \gamma_{\min}$, then (27) cannot hold; if $\Gamma = \gamma_{\min}$, then the encoder can only use inputs in \mathcal{X}' , and the zero-error feedback capacity is thus that of the channel with input alphabet \mathcal{X}' and without a cost constraint; and if $\Gamma \geq \gamma_{\max}$, then (27) always holds, and the cost constraint can be ignored.

As we argue next,

$$C_{f,0}(\Gamma) \geq \frac{\Gamma - \gamma_{\min}}{\gamma_{\max} - \gamma_{\min}} C_{f,0}. \quad (29)$$

In fact $C_{f,0}(\cdot)$ is nondecreasing and concave on $[\gamma_{\min}, \gamma_{\max}]$. Indeed, we can divide the blocklength- n transmission into two frames, Frame 1 and Frame 2, with the former of αn channel uses and the latter of $(1 - \alpha)n$ channel uses, where

$$\alpha = \frac{\gamma_{\max} - \Gamma}{\gamma_{\max} - \gamma_{\min}}.$$

If in Frame 1 the encoder repeatedly transmits an element of \mathcal{X}' , then the cost constraint will be satisfied irrespective of the inputs in Frame 2. Those can thus be chosen to achieve the unconstrained capacity $C_{f,0}$, with the resulting rate being the RHS of (29). This proves (29).

It follows from (29) that $C_{f,0}(\Gamma)$ is positive iff $C_{f,0}$ is positive. By adapting the proof of Theorem 2.4 to account for the cost constraint (27) (see Appendix J), we obtain the following generalization of Theorems 2.3 and 2.4:

Theorem 2.20. *Given any Γ satisfying (28), a necessary and sufficient condition for $C_{f,0}(\Gamma)$ to be positive is (10). If $C_{f,0}(\Gamma)$ is positive, then*

$$C_{f,0}(\Gamma) = \min_{P_S} \max_{\substack{P_{U,X|S}: \\ \mathbb{E}[\gamma(X)] \leq \Gamma}} \min_{P_{Y|U,X,S}: P_{Y|U=u,X,S} \in \mathcal{P}(W), \forall u \in \mathcal{U}} I(U; Y) - I(U; S), \quad (30)$$

where U is an auxiliary chance variable taking values in a finite set \mathcal{U} , the expectation is computed w.r.t the joint PMF $P_S \times P_{U,X|S}$, and the mutual informations are computed w.r.t. the joint PMF $P_S \times P_{U,X|S} \times P_{Y|U,X,S}$. Restricting X to be a function of U and S , i.e., $P_{U,X|S}$ to have the form (12), does not change the RHS of (30), nor does restricting the cardinality of \mathcal{U} to (13).

Specializing Theorem 2.20 to the state-less case, we obtain:

Corollary 2.21. *For a state-less DMC $W(y|x)$ and any Γ satisfying (28), $C_{f,0}(\Gamma)$ is positive iff (2) holds. If $C_{f,0}(\Gamma)$ is positive, then it is given by*

$$C_{f,0}(\Gamma) = \max_{\substack{P_X: \\ \mathbb{E}[\gamma(X)] \leq \Gamma}} \min_{P_{Y|X} \in \mathcal{P}(W)} I(X; Y), \quad (31)$$

where the expectation is computed w.r.t. the PMF P_X and the mutual information w.r.t. the joint PMF $P_X \times P_{Y|X}$.

Proof of Corollary 2.21. This follows from Theorem 2.20 when we consider an SD-DMC $W(y|x, s)$ with a single state, i.e., $|\mathcal{S}| = 1$, whose transition law is

$$W(y|x, s) = W(y|x), \quad \forall (x, s, y) \in \mathcal{X} \times \mathcal{S} \times \mathcal{Y}, \quad (32)$$

because on this channel SI is useless, (10) is equivalent to (2), and the RHS of (30) equals that of (31). \square

The RHS of (31) is a natural generalization of Ahlswede's capacity formula (4) to the setting with the cost constraint (27). Since Ahlswede's capacity formula (4) is an alternative form for Shannon's capacity formula (3), one might wonder whether the RHS of (31) can also be expressed as the "natural" generalization of Shannon's formula (3), namely as

$$\max_{\substack{P_X: \\ \mathbb{E}[\gamma(X)] \leq \Gamma}} \min_{y \in \mathcal{Y}} -\log \sum_{x \in \mathcal{X}: W(y|x) > 0} P_X(x), \quad (33)$$

where the expectation is computed w.r.t. the PMF P_X . The answer is no:

Remark 2.22. *For any $\Gamma \geq \gamma_{\min}$ and every state-less DMC $W(y|x)$*

$$\max_{\substack{P_X: \\ \mathbb{E}[\gamma(X)] \leq \Gamma}} \min_{y \in \mathcal{Y}} -\log \sum_{x \in \mathcal{X}: W(y|x) > 0} P_X(x) \leq \max_{\substack{P_X: \\ \mathbb{E}[\gamma(X)] \leq \Gamma}} \min_{P_{Y|X} \in \mathcal{P}(W)} I(X; Y), \quad (34)$$

where the expectations are computed w.r.t. the PMF P_X and the mutual information w.r.t. the joint PMF $P_X \times P_{Y|X}$. The inequality can be strict.

The inequality (34) is proved in Appendix K. That it can be strict follows from the example below:

Example 2.23. *Suppose*

$$\mathcal{X} = \mathcal{Y} = \{0, 1\};$$

that for every $(x, y) \in \mathcal{X} \times \mathcal{Y}$

$$W(y|x) = \mathbb{1}_{\{y=x\}}, \quad (35)$$

$$\gamma(x) = x; \quad (36)$$

and that $0 < \Gamma < 1/2$. The RHS of (34) evaluates to

$$\max_{\substack{P_X: \\ \mathbb{E}[\gamma(X)] \leq \Gamma}} \min_{P_{Y|X} \in \mathcal{P}(W)} I(X; Y) = h_b(\Gamma); \quad (37)$$

the LHS of (34) evaluates to

$$\max_{\substack{P_X: \\ \mathbb{E}[\gamma(X)] \leq \Gamma}} \min_{y \in \mathcal{Y}} -\log \sum_{x \in \mathcal{X}: W(y|x) > 0} P_X(x) = -\log(1 - \Gamma); \quad (38)$$

and

$$-\log(1 - \Gamma) < h_b(\Gamma), \quad 0 < \Gamma < 1/2. \quad (39)$$

The following may explain why the inequality in (34) can be strict. Recall Shannon's sequential coding scheme [1], which achieves the zero-error feedback capacity (3) of the state-less DMC: The encoder selects some PMF P_X , and, before every channel use, it maps a fraction of approximately $P_X(x)$ of the survivor set to the input symbol x . If the channel output is $y \in \mathcal{Y}$, then the survivor set is reduced by a factor of nearly

$$\left(\sum_{x \in \mathcal{X}: W(y|x) > 0} P_X(x) \right)^{-1}. \quad (40)$$

The generalization (33) of Shannon's capacity formula (3) is obtained when the PMF P_X is restricted to satisfy $\mathbb{E}[\gamma(X)] \leq \Gamma$. As the following argument suggests, a more adaptive coding scheme may be required in the presence of the cost constraint (27). To see why, fix some PMF P_X w.r.t. which $\mathbb{E}[\gamma(X)] \leq \Gamma$, and let $y^* \in \mathcal{Y}$ maximize

$$\sum_{x \in \mathcal{X}: W(y^*|x) > 0} P_X(x). \quad (41)$$

If the cost of every input symbol $x \in \mathcal{X}$ for which $W(y^*|x) > 0$ is smaller than Γ , then the cost constraint loosens for the remaining channel uses, and the encoder should take advantage of this.

2.6 Constrained States

This section provides some insight into how cost constraints on the channel states affect the zero-error feedback capacity of the SD-DMC $W(y|x, s)$ with acausal SI. Consider some nonnegative “cost-function” $\lambda: \mathcal{S} \rightarrow \mathbb{R}_0^+$, define

$$\lambda_{\min} = \min_{s \in \mathcal{S}} \lambda(s) \quad \text{and} \quad \lambda_{\max} = \max_{s \in \mathcal{S}} \lambda(s),$$

and let

$$\lambda_{\min} < \Lambda < \lambda_{\max}. \tag{42}$$

Like the cost constraint (27) on the channel inputs, where we restrict Γ to (28), we restrict Λ to (42), because all other values of Λ are uninteresting. In the following, we shall consider two different cost constraints on the channel states.

The first is that, at every blocklength n , the channel states’ average cost

$$\lambda^{(n)}(S^n) = \frac{1}{n} \sum_{i=1}^n \lambda(S_i)$$

satisfy the cost constraint

$$\lambda^{(n)}(S^n) \leq \Lambda \tag{43}$$

for some given Γ satisfying (42). Let $C_{f,0}^{(1)}(\Lambda)$ denote the zero-error feedback capacity with acausal SI subject to (43). Unlike the cost constraint on the channel inputs (27), the cost constraint on the channel states (43) affects not only the formula for $C_{f,0}$ when it is positive but also whether $C_{f,0}$ is positive. The reason for this is that the time-sharing argument of Section 2.5 does not work for the adversarial state selector: since the state is revealed acausally to the encoder, if the state selector chooses only “benign” states of low cost during Frame 1 and only “hurtful” states of high cost during Frame 2, then the encoder can concentrate its transmission in the first frame, where the state assumes only “benign” realizations of low cost.

Indeed, the cost constraint (43) can increase the zero-error feedback capacity with acausal SI from zero:

Remark 2.24. *Even when Λ satisfies (42), the zero-error feedback capacity of an SD-DMC with acausal SI can be zero in the absence of a state cost-constraint yet be positive in its presence.*

We prove Remark 2.24 by means of the following example:

Example 2.25. Consider a deterministic SD-DMC $W(y|x, s)$ over the binary alphabets $\mathcal{X} = \mathcal{S} = \mathcal{Y} = \{0, 1\}$ with the state cost-function

$$\lambda(s) = s, \quad s \in \mathcal{S}. \quad (44)$$

Let the output corresponding to the input x and the state s be the single element of the set $\mathcal{Y}_{x,s}$ that is given in Table 2

$$\{y \in \mathcal{Y} : W(y|x, s) > 0\} = \mathcal{Y}_{x,s}, \quad \forall (x, s) \in \mathcal{X} \times \mathcal{S}. \quad (45)$$

Since (10) does not hold for this channel, Theorem 2.3 implies that $C_{f,0}$ is zero. However, as shown in Appendix L, $C_{f,0}^{(1)}(\Lambda)$ is positive when $\Lambda > 0$ is sufficiently small so that

$$\Lambda + h_b(\Lambda) < 1. \quad (46)$$

This holds also in the absence of feedback: because $W(y|x, s)$ is $\{0, 1\}$ -valued, the encoder can compute the output from the state (which is revealed to it acausally) and from the input (that it produces), and feedback does not, therefore, increase capacity.

$\mathcal{Y}_{x,s}$		s	
		0	1
x	0	$\{0\}$	$\{1\}$
	1	$\{1\}$	$\{1\}$

Table 2: Nonzero transitions of the SD-DMC in Example 2.25.

If $W(y|x, s)$ satisfies (10), i.e., if $C_{f,0}$ is positive in the absence of a state cost-constraint, then we can adapt the proof of Theorem 2.4 to account for the cost constraint (43) and to thus express $C_{f,0}^{(1)}(\Lambda)$ as the “natural” generalization of (11), i.e., as the RHS of (50) ahead. However, by Remark 2.24 the capacity can be positive also when (10) does not hold; and for this case we do not have a generalization of Theorems 2.3 and 2.4. The difficulty in extending Theorems 2.3 and 2.4 to this case is that the cost constraint (43) allows the adversarial state selector to choose whichever states it likes in βn epochs, where

$$\beta = \frac{\Lambda - \lambda_{\min}}{\lambda_{\max} - \lambda_{\min}}, \quad (47)$$

and these epochs are not revealed to the receiver. This is problematic, because the coding schemes by which we prove the direct parts of Theorems 2.3 and 2.4 comprise multiple short transmission phases. For example, the last block of the coding scheme

by which we prove the direct part of Theorem 2.4 is of negligible length compared to n and consequently also compared to βn , and hence the adversarial state selector is free to choose whichever states it likes during the last block.

The second type of cost constraint we consider is that, for some fixed $l \in \mathbb{N}$ and at every blocklength n , the channel states satisfy the cost constraint

$$\frac{1}{l} \sum_{i=(j-1)l+1}^{jl} \lambda(S_i) \leq \Lambda, \quad (\forall j \in \mathbb{N} \text{ s.t. } jl \leq n). \quad (48)$$

Note that (48) is more stringent than (43), because it constrains the average cost of prespecified l -blocks of consecutive channel states and consequently also the channel states' average over the entire blocklength. The zero-error capacity subject to (48), $C_{f,0}^{(2)}(\Lambda, l)$, depends on l . We define the zero-error feedback capacity of the SD-DMC $W(y|x, s)$ with acausal SI under this type of constraint as

$$\liminf_{l \rightarrow \infty} C_{f,0}^{(2)}(\Lambda, l),$$

and we denote it $C_{f,0}^{(2)}(\Lambda)$. By adapting the proofs of Theorems 2.3 and 2.4 to account for the cost constraint (48) (see Appendix M), we obtain the following single-letter characterization of $C_{f,0}^{(2)}(\Lambda)$:

Theorem 2.26. *Given any Λ satisfying (42), a necessary condition for $C_{f,0}^{(2)}(\Lambda)$ to be positive is that*

$$\left(\forall s, s' \in \mathcal{S} \text{ s.t. } \frac{\lambda(s) + \lambda(s')}{2} \leq \Lambda \right) \quad \exists x, x' \in \mathcal{X} \text{ s.t.} \quad (49)$$

$$\left(W(y|x, s) W(y|x', s') = 0, \quad \forall y \in \mathcal{Y} \right).$$

If this condition holds, then

$$C_{f,0}^{(2)}(\Lambda) = \min_{\substack{P_S: \\ \mathbb{E}[\lambda(S)] \leq \Lambda}} \max_{P_{U,X|S}} \min_{\substack{P_{Y|U,X,S}: \\ P_{Y|U=u,X,S} \in \mathcal{P}(W), \forall u \in \mathcal{U}}} I(U; Y) - I(U; S), \quad (50)$$

where U is an auxiliary chance variable taking values in a finite set \mathcal{U} , the expectation is computed w.r.t the PMF P_S , and the mutual informations are computed w.r.t. the joint PMF $P_S \times P_{U,X|S} \times P_{Y|U,X,S}$. Restricting X to be a function of U and S , i.e., $P_{U,X|S}$ to have the form (12), does not change the RHS of (50), nor does restricting the cardinality of \mathcal{U} to (13).

We do not know whether (49) guarantees that the RHS of (50) be positive, and hence we do not know whether (49) is also sufficient for $C_{f,0}(\Lambda)$ to be positive.

For the deterministic SD-DMC of Example 2.25, Theorem 2.26 yields the following result:

Example 2.27. *For the channel and cost-function of Example 2.25*

$$C_{f,0}^{(2)}(\Lambda) = 1 - \Lambda, \quad 0 \leq \Lambda \leq 1. \quad (51)$$

Proof. Here (49) holds iff $\Lambda < 1$, so the capacity is zero if $\Lambda = 1$. (This could have also been established by noting that the all-one state-sequence results in the output being one irrespective of the input.) If $0 \leq \Lambda < 1$, then the capacity is

$$C_{f,0}^{(2)}(\Lambda) = \min_{\substack{P_S: \\ P_S(1) \leq \Lambda}} \max_{P_{U,X|S}} I(U; Y) - I(U; S) \quad (52)$$

$$= \min_{\substack{P_S: \\ P_S(1) \leq \Lambda}} \sum_{s \in \mathcal{S}} P_S(s) \log |\{y \in \mathcal{Y}: \exists x \in \mathcal{X} \text{ s.t. } W(y|x, s) > 0\}| \quad (53)$$

$$= 1 - \Lambda, \quad (54)$$

where the mutual informations are computed w.r.t. the joint PMF $P_S \times P_{U,X|S} \times W$, and the first two equalities can be proved similarly as in Appendix C. \square

3 Selected Proofs

This section contains the proofs of the results in Section 2.1: Theorem 2.3 is proved in Section 3.1; Theorem 2.4 in Section 3.2; and Theorem 2.7 in Section 3.3.

3.1 A Proof of Theorem 2.3

The proof consists of a direct and a converse part. We first establish the direct part. In fact, we prove the following stronger result:

Remark 3.1. *Consider an SD-DMC $W(y|x, s)$ with feedback whose encoder is furnished with acausal SI. If (10) holds, then n_{bit} channel uses suffice for the error-free transmission of a bit, where n_{bit} is 1 if $|\mathcal{S}| = 1$, and is otherwise upper-bounded by³*

$$\frac{2|\mathcal{Y}| \log |\mathcal{S}| - \log |\mathcal{Y}|}{\log |\mathcal{Y}| - \log (|\mathcal{Y}| - 1)} + 1 + 2|\mathcal{Y}|. \quad (55)$$

³Note that all logarithms in (55) are nonnegative, because (10) implies that $|\mathcal{Y}| \geq 2$.

The direct part of Theorem 2.3 follows from Remark 3.1, because if (10) is satisfied, then, by Remark 3.1,

$$C_{f,0} \geq 1/n_{\text{bit}} > 0. \quad (56)$$

In proving Remark 3.1 we focus on the case $|\mathcal{S}| \geq 2$, because the case $|\mathcal{S}| = 1$ follows directly from Shannon [1]. (In this case (10) is equivalent to (2).)

Before we prove Remark 3.1, we briefly describe the coding scheme that we propose. Because the zero-error capacity of the SD-DMC $W(y|x, s)$ with acausal SI can be zero without feedback but positive with feedback (Theorem 2.7), it is not always possible to transmit a single bit error-free in only one channel use (Corollary 2.8). Our scheme thus requires more than one channel use, and it utilizes the feedback link.

The scheme has two phases. Phase 1 is not used to convey the bit but rather to reduce the decoder's ambiguity about the Phase-2 state-sequence. This is attained with an adaptive feedback code reminiscent of the one used in the first phase of Shannon's coding scheme for the stateless DMC [1]. But in our Phase 1, the encoder utilizes the Phase-1 state-sequence (albeit only causally). After Phase 1 the decoder computes the set of Phase-2 state-sequences of positive posterior probability given the Phase-1 outputs. This set can also be computed by the encoder thanks to the Phase-1 feedback. This enables the encoder to transmit the bit error-free in Phase 2. The feedback link is not used in Phase 2.

The condition in Theorem 2.3 ensures that Phase 1 and 2 are feasible. As we shall see, Phase 1 is feasible iff (14) holds, whereas Phase 2 is feasible iff (10) holds, where by Remarks 2.5 and 2.6

$$(10) \implies (14) \quad \text{and} \quad (10) \not\Leftarrow (14),$$

so feasibility is easier to attain in Phase 1 than in Phase 2.

Proof of Remark 3.1. The case $|\mathcal{S}| = 1$ follows from Shannon [1], and we hence assume that $|\mathcal{S}| \geq 2$. To transmit a single bit $m \in \{0, 1\}$, we divide the blocklength- n_{bit} transmission into Phase 1 and Phase 2 of n_1 and n_2 channel uses, where

$$n_{\text{bit}} = n_1 + n_2. \quad (57)$$

For now, $(n_{\text{bit}}, n_1, n_2)$ could be any triple of positive integers satisfying (57). At the end of the proof, we shall exhibit a choice of the triple for which the transmission is error-free and n_{bit} is upper-bounded by (55). Before we do that, we describe Phase 1 and Phase 2, beginning with Phase 1.

Let $\mathcal{S}^{n_1+n_2}$ denote the set of possible length- $(n_1 + n_2)$ state-sequences, and let \mathcal{S}^{n_2} denote the set of possible state sequences occurring during Phase 2. Before the transmission begins, the encoder observes the entire state sequence $\mathcal{S}^{n_1+n_2}$. The goal of Phase 1 is to produce a random subset $\mathcal{S}_{n_1} \subseteq \mathcal{S}^{n_2}$ with the following three properties: 1) \mathcal{S}_{n_1} is determined by the Phase-1 outputs Y_1, \dots, Y_{n_1} , so both encoder and decoder know \mathcal{S}_{n_1} before Phase 2 begins; 2) with probability one \mathcal{S}_{n_1} contains the Phase-2 state-sequence $\mathcal{S}_{n_1+1}^{n_1+n_2}$; and 3) the cardinality of \mathcal{S}_{n_1} is upper-bound by

$$|\mathcal{S}_{n_1}| \leq \left(\frac{|\mathcal{Y}| - 1}{|\mathcal{Y}|} \right)^{n_1} |\mathcal{S}|^{n_2} + |\mathcal{Y}|. \quad (58)$$

To that end we partition the set $\mathcal{S}_0 = \mathcal{S}^{n_2}$ into $|\mathcal{Y}|$ different subsets whose size is between $\lfloor |\mathcal{S}_0|/|\mathcal{Y}| \rfloor$ and $\lceil |\mathcal{S}_0|/|\mathcal{Y}| \rceil$. We index the $|\mathcal{Y}|$ subsets by the output alphabet \mathcal{Y} and reveal the result to the encoder and decoder. To every pair $(s, y) \in \mathcal{S} \times \mathcal{Y}$ we assign an input $x(s, y) \in \mathcal{X}$ for which

$$W(y|x(s, y), s) = 0. \quad (59)$$

Such an $x(s, y)$ exists, because substituting s for both s and s' in (10) demonstrates that (10) implies that there exists a pair of inputs $x', x'' \in \mathcal{X}$ for which

$$W(y|x', s) W(y|x'', s) = 0, \quad \forall y \in \mathcal{Y}, \quad (60)$$

i.e., for which for every $y \in \mathcal{Y}$ either $W(y|x', s)$ or $W(y|x'', s)$ is zero. We can thus choose $x(s, y)$ to be x' when $W(y|x', s)$ is zero and to be x'' when it is not.⁴ If, thanks to its acausal SI, the encoder knows that the Time-1 state S_1 is s and that $\mathcal{S}_{n_1+1}^{n_1+n_2}$ is in the subset of \mathcal{S}_0 indexed by y , then at Time 1 it transmits $x(s, y)$. This choice guarantees by (59) that, upon observing the Time-1 output Y_1 , the decoder will know that the Phase-2 state-sequence is not an element of the subset of \mathcal{S}_0 indexed by Y_1 , and that it is thus in the \mathcal{S}_0 -complement of this subset, which we denote \mathcal{S}_1 . Note that: 1) both encoder and decoder know \mathcal{S}_1 after Channel-Use 1; 2) \mathcal{S}_1 contains $\mathcal{S}_{n_1+1}^{n_1+n_2}$; and 3) the cardinality of \mathcal{S}_1 is upper-bounded by

$$|\mathcal{S}_1| \leq |\mathcal{S}_0| - \left\lfloor \frac{|\mathcal{S}_0|}{|\mathcal{Y}|} \right\rfloor = \left\lceil \frac{|\mathcal{Y}| - 1}{|\mathcal{Y}|} |\mathcal{S}_0| \right\rceil \leq \frac{|\mathcal{Y}| - 1}{|\mathcal{Y}|} |\mathcal{S}_0| + 1. \quad (61)$$

Phase 1 continues in the same fashion: Let $i \in [2 : n_1]$, and assume that the first $i - 1$ channel uses have produced a random subset \mathcal{S}_{i-1} of \mathcal{S}^{n_2} with the following three properties: 1) both encoder and decoder know \mathcal{S}_{i-1} after Channel-Use $(i - 1)$; 2) \mathcal{S}_{i-1}

⁴This is nothing else but (10) \implies (14), which follows from Remarks 2.5 and 2.6.

contains $S_{n_1+1}^{n_1+n_2}$; and 3) the cardinality of \mathcal{S}_{i-1} is upper-bounded by

$$|\mathcal{S}_{i-1}| \leq \frac{|\mathcal{Y}| - 1}{|\mathcal{Y}|} |\mathcal{S}_{i-2}| + 1. \quad (62)$$

After Channel-Use $(i - 1)$, we partition \mathcal{S}_{i-1} into $|\mathcal{Y}|$ different subsets whose size is between $\lfloor |\mathcal{S}_{i-1}|/|\mathcal{Y}| \rfloor$ and $\lceil |\mathcal{S}_{i-1}|/|\mathcal{Y}| \rceil$. We index the subsets by the elements of the output alphabet \mathcal{Y} and reveal the result to the encoder and decoder. If, thanks to its acausal SI, the encoder knows that the Time- i state S_i is s and that $S_{n_1+1}^{n_1+n_2}$ is an element of the subset of \mathcal{S}_{i-1} indexed by y , then it transmits $x(s, y)$ at Time i . This choice guarantees by (59) that, upon observing the Time- i output Y_i , the decoder will know that the Phase-2 state-sequence is not an element of the subset indexed by Y_i , and that it is thus in the \mathcal{S}_{i-1} -complement of this subset, which we denote \mathcal{S}_i . Note that: 1) both encoder and decoder know \mathcal{S}_i after Channel-Use i ; 2) \mathcal{S}_i contains $S_{n_1+1}^{n_1+n_2}$; and 3) the cardinality of \mathcal{S}_i is upper-bounded by

$$|\mathcal{S}_i| \leq |\mathcal{S}_{i-1}| - \left\lfloor \frac{|\mathcal{S}_{i-1}|}{|\mathcal{Y}|} \right\rfloor = \left\lceil \frac{|\mathcal{Y}| - 1}{|\mathcal{Y}|} |\mathcal{S}_{i-1}| \right\rceil \leq \frac{|\mathcal{Y}| - 1}{|\mathcal{Y}|} |\mathcal{S}_{i-1}| + 1. \quad (63)$$

Since this holds for every $i \in [1 : n_1]$, the goal of Phase 1 is attained, and the first n_1 channel uses produce a random subset \mathcal{S}_{n_1} of \mathcal{S}^{n_2} with the following three properties: 1) both encoder and decoder know \mathcal{S}_{n_1} before Phase 2 begins; 2) \mathcal{S}_{n_1} contains the Phase-2 state-sequence $S_{n_1+1}^{n_1+n_2}$; and 3) the cardinality of \mathcal{S}_{n_1} is upper-bound by

$$|\mathcal{S}_{n_1}| \leq \left(\frac{|\mathcal{Y}| - 1}{|\mathcal{Y}|} \right)^{n_1} |\mathcal{S}_0| + \sum_{i=0}^{n_1-1} \left(\frac{|\mathcal{Y}| - 1}{|\mathcal{Y}|} \right)^i \quad (64)$$

$$= \left(\frac{|\mathcal{Y}| - 1}{|\mathcal{Y}|} \right)^{n_1} |\mathcal{S}|^{n_2} + \frac{|\mathcal{Y}|^{n_1} - (|\mathcal{Y}| - 1)^{n_1}}{|\mathcal{Y}|^{n_1} - (|\mathcal{Y}| - 1)|\mathcal{Y}|^{n_1-1}} \quad (65)$$

$$= \left(\frac{|\mathcal{Y}| - 1}{|\mathcal{Y}|} \right)^{n_1} |\mathcal{S}|^{n_2} + \frac{|\mathcal{Y}|^{n_1} - (|\mathcal{Y}| - 1)^{n_1}}{|\mathcal{Y}|^{n_1-1}} \quad (66)$$

$$\leq \left(\frac{|\mathcal{Y}| - 1}{|\mathcal{Y}|} \right)^{n_1} |\mathcal{S}|^{n_2} + |\mathcal{Y}|. \quad (67)$$

We next turn to Phase 2 whose goal is to transmit the bit error-free. To that end the encoder allocates to every bit value $m \in \{0, 1\}$ and every state sequence \mathbf{s} in \mathcal{S}_{n_1} a length- n_2 codeword $\mathbf{x}(m, \mathbf{s})$, where the codewords are chosen so that

$$\forall \mathbf{s}, \mathbf{s}' \in \mathcal{S}_{n_1} \quad \exists i \in [1 : n_2] \text{ s.t. } \left(W(y|x_i(0, \mathbf{s}), s_i) W(y|x_i(1, \mathbf{s}'), s'_i) = 0, \forall y \in \mathcal{Y} \right). \quad (68)$$

(We will shortly show how this can be done.) If the value of the bit to be sent is $m \in \{0, 1\}$ and if the Phase-2 state-sequence is \mathbf{s} , then the encoder transmits in Phase 2 the

codeword $\mathbf{x}(m, \mathbf{s})$. Condition (68) implies that, upon observing the realization $\mathbf{y} \in \mathcal{Y}^{n_2}$ of the Phase-2 output-sequence $Y_{n_1+1}^{n_1+n_2}$, the decoder, who knows \mathcal{S}_{n_1} and the codewords $\{\mathbf{x}(\tilde{m}, \tilde{\mathbf{s}})\}$, can determine the value of m error-free, because for the true realization $\mathbf{s} \in \mathcal{S}_{n_1}$ of the Phase-2 state-sequence

$$\prod_{i=1}^{n_2} W(y_i | x_i(m, \mathbf{s}), s_i) > 0, \quad (69)$$

whereas (68) implies for $m' \neq m$

$$\prod_{i=1}^{n_2} W(y_i | x_i(m', \tilde{\mathbf{s}}), \tilde{s}_i) = 0, \quad \forall \tilde{\mathbf{s}} \in \mathcal{S}_{n_1}. \quad (70)$$

The decoder can thus calculate $\prod_i W(y_i | x_i(\tilde{m}, \tilde{\mathbf{s}}), \tilde{s}_i)$ for each $\tilde{m} \in \{0, 1\}$ and $\tilde{\mathbf{s}} \in \mathcal{S}_{n_1}$ and produce the \tilde{m} for which this product is positive for some $\tilde{\mathbf{s}} \in \mathcal{S}_{n_1}$.

One (inefficient) way to achieve (68) is the following. Let x^* be an arbitrary fixed element of \mathcal{X} , and for every pair $s, s' \in \mathcal{S}$ choose a pair $x(s, s'), x'(s, s') \in \mathcal{X}$ for which

$$W(y | x(s, s'), s) W(y | x'(s, s'), s') = 0, \quad \forall y \in \mathcal{Y}. \quad (71)$$

By (10) such a pair $x(s, s'), x'(s, s')$ exists. Now choose

$$n_2 \geq |\mathcal{S}_{n_1}|^2; \quad (72)$$

allocate to every ordered pair $(\mathbf{s}, \mathbf{s}') \in \mathcal{S}_{n_1} \times \mathcal{S}_{n_1}$ a different index $i \in [1 : |\mathcal{S}_{n_1}|^2]$; and for the allocated index i choose $x_i(0, \mathbf{s}) = x(s_i, s'_i)$ and $x_i(1, \mathbf{s}') = x'(s_i, s'_i)$, and thus guarantee, by (71), that

$$\left(W(y | x_i(0, \mathbf{s}), s_i) W(y | x_i(1, \mathbf{s}'), s'_i) = 0, \quad \forall y \in \mathcal{Y} \right). \quad (73)$$

The above specifies $|\mathcal{S}_{n_1}|$ out of $n_2 \geq |\mathcal{S}_{n_1}|^2$ symbols of each codeword $\mathbf{x}(m, \mathbf{s})$. How we choose the other $n_2 - |\mathcal{S}_{n_1}|$ symbols is immaterial. To be explicit, we choose each of them to be x^* . The described choice of the codewords $\{\mathbf{x}(m, \mathbf{s})\}$ clearly satisfies (68). Hence, it would only remain to exhibit some choice of the triple $(n_{\text{bit}}, n_1, n_2)$ satisfying (57) and (72). This can be done using (58), but the resulting value of n_{bit} need not be upper-bounded by (55). To fix this, we allocate the indices more efficiently. Note that for every $i \in [1 : |\mathcal{S}_{n_1}|^2]$ the above choice of the codewords $\{\mathbf{x}(m, \mathbf{s})\}$ allocates meaningful values to the i -th symbols of only two codewords, namely $\mathbf{x}(0, \mathbf{s})$ and $\mathbf{x}(1, \mathbf{s}')$, where $(\mathbf{s}, \mathbf{s}')$ is the ordered pair to which Index i has been allocated. More efficiently, we can allocate the same index i to several distinct pairs $(\mathbf{s}, \mathbf{s}')$. (Still, we let $x_i(0, \mathbf{s}) = x(s_i, s'_i)$ and $x_i(1, \mathbf{s}') = x'(s_i, s'_i)$ when Index i has been allocated to the ordered pair $(\mathbf{s}, \mathbf{s}')$, and we

choose each codeword symbol that has not been assigned a value to be x^* .) This works whenever any two distinct pairs $(\mathbf{s}, \mathbf{s}')$, $(\tilde{\mathbf{s}}, \tilde{\mathbf{s}}')$ that are allocated the same index i satisfy $\mathbf{s} \neq \tilde{\mathbf{s}}$ and $\mathbf{s}' \neq \tilde{\mathbf{s}}'$, because then every codeword symbol $x_i(m, \mathbf{s})$ is assigned exactly one value. An efficient way to allocate the indices and guarantee that this requirement is met is the following. Instead of (72), choose any integer n_2 that satisfies

$$n_2 \geq |\mathcal{S}_{n_1}|. \quad (74)$$

(An explicit choice for which n_{bit} is upper-bounded by (55) will be given in (77).) Index the elements of \mathcal{S}_{n_1} by $[1 : |\mathcal{S}_{n_1}|]$, where $\mathbf{s}(j)$ denotes the element of \mathcal{S}_{n_1} indexed by j . Allocate to every ordered pair $(\mathbf{s}(k), \mathbf{s}(\ell))$, where $k, \ell \in [1 : |\mathcal{S}_{n_1}|]$, the index

$$i(k, \ell) = (\ell - k \bmod |\mathcal{S}_{n_1}|) + 1, \quad (75)$$

which clearly satisfies

$$i \in [1 : |\mathcal{S}_{n_1}|] \subseteq [1 : n_2]. \quad (76)$$

By (75) any two distinct pairs $(\mathbf{s}(k), \mathbf{s}(\ell))$, $(\mathbf{s}(k'), \mathbf{s}(\ell'))$ that are allocated the same index i satisfy $k \neq k'$ and $\ell \neq \ell'$, so $\mathbf{s}(k) \neq \mathbf{s}(k')$ and $\mathbf{s}(\ell) \neq \mathbf{s}(\ell')$.

To conclude the direct part, it remains to exhibit some choice of the triple $(n_{\text{bit}}, n_1, n_2)$ satisfying (57) and (74). By (58) these are satisfied if

$$n_1 = \left\lceil \frac{2|\mathcal{Y}|\log|\mathcal{S}| - \log|\mathcal{Y}|}{\log|\mathcal{Y}| - \log(|\mathcal{Y}| - 1)} \right\rceil, \quad (77a)$$

$$n_2 = 2|\mathcal{Y}|, \quad (77b)$$

$$n_{\text{bit}} = \left\lceil \frac{2|\mathcal{Y}|\log|\mathcal{S}| - \log|\mathcal{Y}|}{\log|\mathcal{Y}| - \log(|\mathcal{Y}| - 1)} \right\rceil + 2|\mathcal{Y}|, \quad (77c)$$

and for this choice n_{bit} is upper-bounded by (55). \square

We next prove the converse part of Theorem 2.3.

Converse Part. To show that (10) is necessary for $C_{f,0}$ to be positive, we need to prove that if (10) does not hold, i.e., if there exists a pair of states $s, s' \in \mathcal{S}$ such that

$$\nexists x, x' \in \mathcal{X} \text{ s.t. } (W(y|x, s)W(y|x', s') = 0, \forall y \in \mathcal{Y}), \quad (78)$$

then it is impossible to transmit a single bit error-free. Condition (78) can be alternatively expressed as

$$\forall x, x' \in \mathcal{X} \quad \exists y \in \mathcal{Y}: W(y|x, s)W(y|x', s') > 0, \quad (79)$$

which makes the claim almost obvious. Indeed, (79) implies that, if the state sequence is all s or all s' , then—during every channel use and irrespective of the inputs x, x' that we choose—the pairs (x, s) and (x', s') can produce the same output. This implies that for every pair of messages $m, m' \in \mathcal{M}$ and every encoding mappings there exists an output sequence of positive probability conditional on each of the following two events: 1) the message is m , and the state sequence is all s ; or 2) the message is m' , and the state sequence is all s' .

To prove this formally, let the bit take values in the set $\mathcal{M} = \{0, 1\}$, and fix a blocklength n and n encoding mappings

$$f_i: \mathcal{M} \times \mathcal{S}^n \times \mathcal{Y}^{i-1} \rightarrow \mathcal{X}, \quad i \in [1 : n].$$

Denote by $\mathbf{s} \in \mathcal{S}^n$ the all- s and by $\mathbf{s}' \in \mathcal{S}^n$ the all- s' state-sequence, so

$$s_i = s \quad \text{and} \quad s'_i = s', \quad \forall i \in [1 : n]. \quad (80)$$

To show that the mappings do not achieve error-free transmission, we will exhibit an output sequence $\mathbf{y} \in \mathcal{Y}^n$ that for every $i \in [1 : n]$ satisfies

$$W(y_i | f_i(0, \mathbf{s}, y^{i-1}), s_i) W(y_i | f_i(1, \mathbf{s}', y^{i-1}), s'_i) > 0. \quad (81)$$

This will rule out error-free transmission, because if the state sequence is either \mathbf{s} or \mathbf{s}' , then the decoder, not knowing which, cannot recover the bit.

Our construction of $\mathbf{y} \in \mathcal{Y}^n$ is inductive, i.e., we first exhibit a Time-1 output $y_1 \in \mathcal{Y}$ that satisfies (81) for $i = 1$, and we then repeatedly increment i by one (until it reaches n) and exhibit a Time- i output $y_i \in \mathcal{Y}$ that—together with the previously constructed $\{y_j\}_{j \in [1:i-1]}$ —satisfies (81).

We start by exhibiting a Time-1 output $y_1 \in \mathcal{Y}$ that satisfies (81) for $i = 1$. To this end we observe from (79) and (80) that

$$\exists y \in \mathcal{Y} \text{ s.t. } W(y | f_1(0, \mathbf{s}), s_1) W(y | f_1(1, \mathbf{s}'), s'_1) > 0. \quad (82)$$

If y is as promised in (82), then we choose $y_1 = y$ with the result that (81) holds for $i = 1$.

For the inductive step, suppose $\ell \in [2 : n]$, and that we have already constructed $\{y_i\}_{i \in [1:\ell-1]}$ for which (81) holds for every $i \in [1 : \ell - 1]$. We construct a Time- ℓ output $y_\ell \in \mathcal{Y}$ that—together with the previously constructed $\{y_i\}_{i \in [1:\ell-1]}$ —satisfies (81) when we substitute ℓ for i in (81), i.e., we show that

$$\exists y_\ell \in \mathcal{Y} \text{ s.t. } W(y_\ell | f_\ell(0, \mathbf{s}, y^{\ell-1}), s_\ell) W(y_\ell | f_\ell(1, \mathbf{s}', y^{\ell-1}), s'_\ell) > 0. \quad (83)$$

In fact, (83) follows from (79) and (80).

Since the construction goes through for every $\ell \in [1 : n]$, when ℓ reaches n we have constructed an output sequence $\mathbf{y} \in \mathcal{Y}^n$ that for every $i \in [1 : n]$ satisfies (81). \square

3.2 A Proof of Theorem 2.4

As we prove in Appendix D, restricting X to be a function of U and S , i.e., $P_{U,X|S}$ to have the form (12), does not change the RHS of (11), nor does restricting the cardinality of \mathcal{U} to (13) (Lemma D.1). To prove Theorem 2.4 it thus suffices to establish a direct part for the case where the cardinality of \mathcal{U} is restricted to (13) and a converse part for the case where \mathcal{U} is any finite set. We first establish the direct part.

Direct Part. Our coding scheme can be roughly described as follows. We partition the blocklength- n transmission into $B + 1$ blocks, with each of the first B blocks being of length k . Each of these blocks is guaranteed to reduce the “survivor set”—i.e., the set of messages of positive posterior probability given the channel outputs—by at least a factor of nearly

$$\min_{P_S} \max_{P_{U,X|S}} \min_{\substack{P_{Y|U,X,S}: \\ P_{Y|U=u,X,S} \in \mathcal{P}(W), \forall u \in \mathcal{U}}} 2^{k(I(U;Y) - I(U;S))},$$

where U is an auxiliary chance variable taking values in a finite set \mathcal{U} , and where the mutual informations are computed w.r.t. the joint PMF $P_S \times P_{U,X|S} \times P_{Y|U,X,S}$. The parameter B is chosen so that the post-Block- B survivor-set be “small.” The last block further reduces the survivor-set from a small set to a singleton containing the transmitted message. The coding scheme asymptotically achieves the rate on the RHS of (11), because, when B and k are large, the last block is of negligible length compared to Bk and therefore does not affect the code’s asymptotic rate.

In the first B blocks our scheme draws on Dueck’s scheme for zero-error communication over the multiple-access channel with feedback [4]. Dueck’s scheme in turn draws on Ahlswede’s work [3, 5, 6], which was originally motivated by the AVC with feedback, and which on the (state-less) DMC $W(y|x)$ achieves the zero-error feedback capacity (4) [5]. We next describe Blocks 1 through B of Ahlswede’s scheme and then show how to adapt them to the present setting.

Fix positive integers B, k and a k -type P_X on \mathcal{X} . Let $\mathcal{M}_0 \triangleq \mathcal{M}$ be the set of possible messages, and for every $b \in [1 : B]$ let \mathcal{M}_b be the post-Block- b survivor-set, i.e., the (random) set of messages of positive posterior probability given the channel outputs Y^{bk} during the first b blocks. Thus, \mathcal{M}_b is the (random) subset of \mathcal{M}_{b-1}

comprising the messages in \mathcal{M}_{b-1} of positive posterior probability given the Block- b outputs $\mathbf{y}^{(b)} \triangleq Y_{(b-1)k+1}^{bk}$. Ahlswede's scheme is designed so as to guarantee that

$$|\mathcal{M}_b| \lesssim \left(\max_{P_{Y|X} \in \mathcal{P}(W)} 2^{-kI(X;Y)} \right) |\mathcal{M}_{b-1}|, \quad (84)$$

where the mutual information is computed w.r.t. the joint PMF $P_X \times P_{Y|X}$.

For every $b \in [1 : B]$ Ahlswede's Block- b transmission can be described as follows. Thanks to the feedback link, the set \mathcal{M}_{b-1} can be computed by both transmitter and receiver after Block $(b-1)$. They can thus agree on a partition of \mathcal{M}_{b-1} into $|\mathcal{T}_{P_X}^{(k)}|$ message sets whose size is between $\lfloor |\mathcal{M}_{b-1}| / |\mathcal{T}_{P_X}^{(k)}| \rfloor$ and $\lceil |\mathcal{M}_{b-1}| / |\mathcal{T}_{P_X}^{(k)}| \rceil$, and they can agree on a way to associate with each message set a different k -tuple from $\mathcal{T}_{P_X}^{(k)}$. To transmit Message $m \in \mathcal{M}_{b-1}$, the encoder transmits the k -tuple $\mathbf{x}^{(b)} \in \mathcal{T}_{P_X}^{(k)}$ associated with the message set containing m . Based on the Block- b outputs $\mathbf{y}^{(b)}$, the encoder and decoder compute \mathcal{M}_b as follows: they identify all the k -tuples in $\mathcal{T}_{P_X}^{(k)}$ that could have produced the Block- b outputs $\mathbf{y}^{(b)}$, and they compute \mathcal{M}_b as the union of the message sets with which these k -tuples are associated.

We next establish (84), or more precisely that

$$|\mathcal{M}_b| \leq \left(\max_{P_{Y|X} \in \mathcal{P}(W)} 2^{-k(I(X;Y) - \alpha_k)} \right) |\mathcal{M}_{b-1}|, \quad (85a)$$

whenever

$$|\mathcal{M}_{b-1}| \geq |\mathcal{T}_{P_X}^{(k)}|, \quad (85b)$$

where the mutual information is computed w.r.t. the joint PMF $P_X \times P_{Y|X}$, and where α_k is given by

$$\alpha_k = \frac{\log(1+k)|\mathcal{X}|(1+|\mathcal{Y}|) + 1}{k} \quad (86)$$

and hence converges to zero as k tends to infinity. To this end assume that (85b) holds and note that, with probability one, the empirical type of the pair of Block- b inputs and outputs $(\mathbf{x}^{(b)}, \mathbf{y}^{(b)})$ satisfies

$$P_{\mathbf{x}^{(b)}} = P_X, \quad (87a)$$

$$\left(W(y|x) = 0 \right) \implies \left(P_{\mathbf{x}^{(b)}, \mathbf{y}^{(b)}}(x, y) = 0 \right). \quad (87b)$$

This allows us to upper-bound the number of k -tuples in $\mathcal{T}_{P_X}^{(k)}$ that could have produced the observed Block- b outputs $\mathbf{y}^{(b)}$: For every fixed k -type $P_{X,Y}$ on $\mathcal{X} \times \mathcal{Y}$, the number of k -tuples \mathbf{x} that satisfy $(\mathbf{x}, \mathbf{y}^{(b)}) \in \mathcal{T}_{P_{X,Y}}^{(k)}$ cannot exceed $2^{kH(X|Y)}$, where the conditional

entropy is computed w.r.t. the joint PMF $P_{X,Y}$ [8, Lemma 2.5]. This, combined with (87) and the fact that the number of k -types on $\mathcal{X} \times \mathcal{Y}$ cannot exceed $(1+k)^{|\mathcal{X}||\mathcal{Y}|}$, implies that the number of k -tuples in $\mathcal{T}_{P_X}^{(k)}$ that could have produced the observed Block- b outputs $\mathbf{y}^{(b)}$ is upper-bounded by

$$2^{\log(1+k)|\mathcal{X}||\mathcal{Y}|} \max_{P_{Y|X} \in \mathcal{P}(W)} 2^{kH(X|Y)}, \quad (88)$$

where the conditional entropy is computed w.r.t. the joint PMF $P_X \times P_{Y|X}$. Every k -tuple from $\mathcal{T}_{P_X}^{(k)}$ is associated with a message set whose size is at most $\lceil |\mathcal{M}_{b-1}| / |\mathcal{T}_{P_X}^{(k)}| \rceil$; and, by the assumption that (85b) holds,

$$\lceil |\mathcal{M}_{b-1}| / |\mathcal{T}_{P_X}^{(k)}| \rceil \stackrel{(a)}{\leq} 2 |\mathcal{M}_{b-1}| / |\mathcal{T}_{P_X}^{(k)}| \quad (89)$$

$$\stackrel{(b)}{\leq} |\mathcal{M}_{b-1}| 2^{-kH(X) + \log(1+k)|\mathcal{X}| + 1}, \quad (90)$$

where (a) follows from (85b); and (b) follows from the inequality $|\mathcal{T}_{P_X}^{(k)}| \geq (1+k)^{-|\mathcal{X}|} 2^{kH(X)}$, where the entropy is computed w.r.t. P_X [8, Lemma 2.3]. From (88) and (90) we obtain (85).

We next sketch our adaption of Ahlswede's scheme to the present setting. For every $b \in [1 : B]$ the Block- b transmission can be described as follows. Before the transmission begins, the encoder is revealed the realization $\mathbf{s}^{(b)} \triangleq S_{(b-1)k+1}^{bk}$ of the Block- b state-sequence. Assume for now that the decoder—while incognizant of $\mathbf{s}^{(b)}$ —knows its empirical type $P_{\mathbf{s}^{(b)}}$: the latter will be conveyed to the decoder error-free in Block $B+1$. Let $\mathcal{M}_0 \triangleq \mathcal{M}$ be the set of possible messages, and let \mathcal{M}_b be the post-Block- b survivor-set, i.e., the (random) subset of \mathcal{M}_{b-1} comprising the messages in \mathcal{M}_{b-1} of positive posterior probability given the Block- b outputs $\mathbf{y}^{(b)}$ and the empirical type $P_{\mathbf{s}^{(b)}}$. Choose some k -type $P_{U,X,S}^{(b)}$ whose \mathcal{S} -marginal $P_S^{(b)}$ equals $P_{\mathbf{s}^{(b)}}$. In the following, unless otherwise specified, all entropies and mutual informations are computed w.r.t. the joint PMF $P_{U,X,S}^{(b)}$. Unlike Ahlswede's Block b , which partitions \mathcal{M}_{b-1} into $|\mathcal{T}_{P_X}^{(k)}|$ message sets and associates with each a different k -tuple from $\mathcal{T}_{P_X}^{(k)}$, we fix some $\epsilon > 0$ and partition \mathcal{M}_{b-1} into

$$\Theta \triangleq \left\lceil 2^{k(H(U|S) - \epsilon)} \right\rceil \quad (91)$$

message sets whose size is between $\lfloor |\mathcal{M}_{b-1}| / \Theta \rfloor$ and $\lceil |\mathcal{M}_{b-1}| / \Theta \rceil$; and we associate with each message set a different bin from the bins

$$\mathcal{B}_\ell \subseteq \mathcal{T}_{P_U^{(b)}}^{(k)}, \quad \ell \in [1 : \Theta],$$

where the bins $\{\mathcal{B}_\ell\}_{\ell \in [1:\Theta]}$ are pairwise disjoint subsets of $\mathcal{T}_{P_U^{(b)}}^{(k)}$

$$\mathcal{B}_\ell \cap \mathcal{B}_{\ell'} = \emptyset, \quad (\forall \ell, \ell' \in [1:\Theta] \text{ s.t. } \ell' \neq \ell), \quad (92a)$$

and where each bin “covers” $\mathcal{T}_{P_S^{(b)}}^{(k)}$ in the sense that

$$\forall (\mathbf{s}, \ell) \in \mathcal{T}_{P_S^{(b)}}^{(k)} \times [1:\Theta] \quad \exists \mathbf{u} \in \mathcal{B}_\ell \text{ s.t. } (\mathbf{u}, \mathbf{s}) \in \mathcal{T}_{P_{U,S}^{(b)}}^{(k)}. \quad (92b)$$

(Lemma 3.2 ahead guarantees the existence of such bins whenever k is sufficiently large.) To transmit Message $m \in \mathcal{M}_{b-1}$, the encoder picks from the bin that is associated with the message set containing m a k -tuple $\mathbf{u}^{(b)}$ satisfying $(\mathbf{u}^{(b)}, \mathbf{s}^{(b)}) \in \mathcal{T}_{P_{U,S}^{(b)}}^{(k)}$. (By (92b) such a k -tuple $\mathbf{u}^{(b)}$ exists.) It then chooses as the Block- b channel-inputs some k -tuple $\mathbf{x}^{(b)}$ satisfying $(\mathbf{u}^{(b)}, \mathbf{x}^{(b)}, \mathbf{s}^{(b)}) \in \mathcal{T}_{P_{U,X,S}^{(b)}}^{(k)}$. (This is possible, because $\mathcal{T}_{P_{U,X,S}^{(b)}}^{(k)}$ is not empty since $P_{U,X,S}^{(b)}$ is a k -type, and because, by (92b), $(\mathbf{u}^{(b)}, \mathbf{s}^{(b)}) \in \mathcal{T}_{P_{U,S}^{(b)}}^{(k)}$.) Based on the Block- b outputs $\mathbf{y}^{(b)}$ and the empirical type $P_{\mathbf{s}^{(b)}}$, the encoder and decoder compute \mathcal{M}_b as follows. First, they identify all the k -tuples in $\mathcal{T}_{P_U^{(b)}}^{(k)}$ that could have produced the observed Block- b outputs $\mathbf{y}^{(b)}$. Then, they determine all the bins that contain at least one of the identified k -tuples. Finally, they compute \mathcal{M}_b as the union of the message sets with which these bins are associated.⁵

Using arguments similar to those for the state-less DMC, we next show that

$$|\mathcal{M}_b| \leq \left(\max_{\substack{P_{Y|U,X,S}: \\ P_{Y|U=u,X,S} \in \mathcal{P}(W), \forall u \in \mathcal{U}}} 2^{-k(I(U;Y) - I(U;S) - (\epsilon + \beta_k))} \right) |\mathcal{M}_{b-1}|, \quad (93a)$$

whenever

$$|\mathcal{M}_{b-1}| \geq 2^{k(H(U|S) - \epsilon)}, \quad (93b)$$

⁵Our Blocks 1 through B are reminiscent of Merhav and Weissman’s ϵ -error scheme for the state-dependent DMC with acausal SI and feedback to the encoder [3, Section III], which also draws on [5, 6]. Unlike the ϵ -error scheme, our Block b must, however, reduce \mathcal{M}_{b-1} with probability one and hence differs from Block b of the ϵ -error scheme in the following three aspects: 1) it can deal with every possible Block- b state-sequence, regardless of whether or not its empirical type is close to the PMF Q of the state; 2) for every fixed k -type $P_{U,S}^{(b)}$ on $\mathcal{U} \times \mathcal{S}$, every Block- b state-sequence $\mathbf{s}^{(b)}$ of empirical type $P_S^{(b)}$, and every message m in \mathcal{M}_{b-1} , the bin allocated to the message set containing m contains some k -tuple $\mathbf{u}^{(b)}$ that satisfies $(\mathbf{u}^{(b)}, \mathbf{s}^{(b)}) \in \mathcal{T}_{P_{U,S}^{(b)}}^{(k)}$; and 3) our Block b can deal with every possible Block- b output-sequence, regardless of whether or not the sequence is typical according to $W(y|x, s)$.

where the mutual informations are computed w.r.t. the joint PMF $P_{U,X,S}^{(b)} \times P_{Y|U,X,S}$, and where β_k is given by

$$\beta_k = \frac{\log(1+k)|\mathcal{U}||\mathcal{Y}| + 1}{k} \quad (94)$$

and hence converges to zero as k tends to infinity. To this end assume that (93b) holds and note that, with probability one, the empirical type of the tuple $(\mathbf{u}^{(b)}, \mathbf{x}^{(b)}, \mathbf{s}^{(b)}, \mathbf{y}^{(b)})$ satisfies

$$P_{\mathbf{u}^{(b)}, \mathbf{x}^{(b)}, \mathbf{s}^{(b)}} = P_{U,X,S}^{(b)}, \quad (95a)$$

$$\left(W(y|x, s) = 0\right) \implies \left(P_{\mathbf{u}^{(b)}, \mathbf{x}^{(b)}, \mathbf{s}^{(b)}, \mathbf{y}^{(b)}}(u, x, s, y) = 0, \forall u \in \mathcal{U}\right). \quad (95b)$$

This allows us to upper-bound the number of k -tuples in $\mathcal{T}_{P_U^{(b)}}^{(k)}$ that could have produced the observed Block- b outputs $\mathbf{y}^{(b)}$: For every fixed k -type $P_{U,Y}$ on $\mathcal{U} \times \mathcal{Y}$, the number of k -tuples \mathbf{u} that satisfy $(\mathbf{u}, \mathbf{y}^{(b)}) \in \mathcal{T}_{P_{U,Y}}^{(k)}$ cannot exceed $2^{kH(U|Y)}$, where the conditional entropy is computed w.r.t. the joint PMF $P_{U,Y}$ [8, Lemma 2.5]. This, combined with (95) and the fact that the number of k -types on $\mathcal{U} \times \mathcal{Y}$ cannot exceed $(1+k)^{|\mathcal{U}||\mathcal{Y}|}$, implies that the number of k -tuples in $\mathcal{T}_{P_U^{(b)}}^{(k)}$ that could have produced the observed Block- b outputs $\mathbf{y}^{(b)}$ is upper-bounded by

$$2^{\log(1+k)|\mathcal{U}||\mathcal{Y}|} \max_{\substack{P_{Y|U,X,S}: \\ P_{Y|U=u,X,S} \in \mathcal{P}(W), \forall u \in \mathcal{U}}} 2^{kH(U|Y)}, \quad (96)$$

where the conditional entropy is computed w.r.t. the joint PMF $P_{U,X,S}^{(b)} \times P_{Y|U,X,S}$. Since the bins are pairwise disjoint (92a), no k -tuple is contained in more than one bin, and (96) is thus also an upper bound on the number of bins that contain at least one k -tuple that could have produced the observed Block- b outputs. Every bin is associated with a message set whose size is at most $\lceil |\mathcal{M}_{b-1}|/\Theta \rceil$; and, by (91) and the assumption that (93b) holds,

$$\lceil |\mathcal{M}_{b-1}|/\Theta \rceil \leq \left\lceil 2^{-k(H(U|S)-\epsilon)} |\mathcal{M}_{b-1}| \right\rceil \leq 2^{-k(H(U|S)-\epsilon)+1} |\mathcal{M}_{b-1}|. \quad (97)$$

From (96), (97), and the fact that

$$H(U|S) - H(U|Y) = I(U; Y) - I(U; S) \quad (98)$$

we obtain (93).

Since $H(U|S) \leq \log |\mathcal{U}|$ and $\epsilon > 0$, it follows from (93) that

$$|\mathcal{M}_b| \leq \left(\max_{\substack{P_{Y|U,X,S}: \\ P_{Y|U=u,X,S} \in \mathcal{P}(W), \forall u \in \mathcal{U}}} 2^{-k(I(U;Y)-I(U;S)-(\epsilon+\beta_k))} \right) |\mathcal{M}_{b-1}|, \quad (99a)$$

whenever

$$|\mathcal{M}_{b-1}| \geq 2^{k \log |\mathcal{U}|}, \quad (99b)$$

where the mutual informations are computed w.r.t. the joint PMF $P_{U,X,S}^{(b)} \times P_{Y|U,X,S}$, and where β_k is defined in (94).

From (99), which holds for every $b \in [1 : B]$, we infer that we can choose B to be the smallest integer for which

$$|\mathcal{M}_B| \leq 2^{k \log |\mathcal{U}|}. \quad (100)$$

In Block $(B + 1)$ we resolve the post-Block- B survivor-set \mathcal{M}_B , and we transmit the empirical types $P_{\mathbf{s}(1)}, \dots, P_{\mathbf{s}(B)}$ of the state sequences pertaining to Blocks 1 through B . It follows from (100) that, when B is large, the number of bits that are needed to resolve \mathcal{M}_B is negligible compared to Bk . Moreover, when k is large, $B \log(1 + k)|\mathcal{S}|$, which upper-bounds the number of bits needed to represent $P_{\mathbf{s}(1)}, \dots, P_{\mathbf{s}(B)}$, is small compared to Bk . If we thus choose B and k sufficiently large, then—compared to Bk —the encoder will only need to transmit few bits error-free in Block $(B + 1)$, and by Remark 3.1 this can be achieved with the length of the last block negligible compared to Bk .

We next describe and analyze our coding scheme in detail, beginning with Blocks 1 through B and ending with the last block. Throughout, we assume that $C_{f,0}$ is positive, which (by Theorem 2.3) is equivalent to the assumption that (10) holds.

For Blocks 1 through B we only provide the missing details. Fix positive integers B, k , some finite set \mathcal{U} of cardinality

$$|\mathcal{U}| \leq |\mathcal{X}|^{|\mathcal{S}|}, \quad (101)$$

and some $\epsilon > 0$. Assume for now that the decoder knows the empirical types $\{P_{\mathbf{s}(b)}\}_{b \in [1:B]}$ of the state sequences $\{\mathbf{s}^{(b)}\}_{b \in [1:B]}$: those will be conveyed to the decoder error-free in Block $B + 1$. Let $\mathcal{M}_0 \triangleq \mathcal{M}$ be the set of possible messages, and for every $b \in [1 : B]$ let \mathcal{M}_b be the post-Block- b survivor-set, i.e., the (random) set of messages of positive posterior probability given the channel outputs Y^{bk} and the empirical types $\{P_{\mathbf{s}(b')}\}_{b' \in [1:b]}$. Thus, \mathcal{M}_b is the subset of \mathcal{M}_{b-1} comprising the messages in \mathcal{M}_{b-1} of positive posterior probability given the Block- b outputs $\mathbf{y}^{(b)}$ and the empirical type $P_{\mathbf{s}(b)}$ of the Block- b state-sequence $\mathbf{s}^{(b)}$. We already described the Block- b transmission for every $b \in [1 : B]$; it only remains to show that we can find bins

$$\mathcal{B}_\ell \subseteq \mathcal{T}_{P_U^{(b)}}^{(k)}, \quad \ell \in [1 : \Theta]$$

such that (92) holds. This follows from the following lemma:

Lemma 3.2. *Let \mathcal{U} and \mathcal{S} be finite sets. For every $\epsilon > 0$ we can find a positive integer $\eta_0 = \eta_0(|\mathcal{U}|, |\mathcal{S}|, \epsilon)$ that will guarantee that, for every $k \geq \eta_0$ and every k -type $P_{U,S}$, there exist a partition $\{\mathcal{B}_\ell\}_{\ell \in [1:\Theta]}$ of the type class $\mathcal{T}_{P_U}^{(k)}$ with the property that*

$$\forall (\mathbf{s}, \ell) \in \mathcal{T}_{P_S}^{(k)} \times [1:\Theta] \quad \exists \mathbf{u} \in \mathcal{B}_\ell \text{ s.t. } (\mathbf{u}, \mathbf{s}) \in \mathcal{T}_{P_{U,S}}^{(k)}, \quad (102)$$

where $\Theta = \lceil 2^{k(H(U|S) - \epsilon)} \rceil$ with $H(U|S)$ being computed w.r.t. the joint PMF $P_{U,S}$.

Proof. See Appendix E. □

By Lemma 3.2 and (101) we can find a positive integer $\eta_0 = \eta_0(|\mathcal{X}|, |\mathcal{S}|, \epsilon)$ that guarantees that, for every $k \geq \eta_0$ and k -type $P_U^{(b)}$, there exist bins

$$\mathcal{B}_\ell \subseteq \mathcal{T}_{P_U^{(b)}}^{(k)}, \quad \ell \in [1:\Theta]$$

satisfying (92).

Henceforth, assume that $k \geq \eta_0$ and that the bins are as above. We next conclude the analysis of Blocks 1 through B by showing that each of these blocks can reduce the survivor set by at least a factor of nearly

$$\min_{P_S} \max_{P_{U,X|S}} \min_{\substack{P_{Y|U,X,S}: \\ P_{Y|U=u,X,S} \in \mathcal{P}(W), \forall u \in \mathcal{U}}} 2^{k(I(U;Y) - I(U;S))},$$

where the mutual informations are computed w.r.t. the joint PMF $P_S \times P_{U,X|S} \times P_{Y|U,X,S}$. To that end recall that if (99b) holds, then $|\mathcal{M}_b|$ can be upper-bounded in terms of $|\mathcal{M}_{b-1}|$ using (99a), where the mutual informations are computed w.r.t. the joint PMF $P_{U,X,S}^{(b)} \times P_{Y|U,X,S}$, and where β_k is defined in (94). Since we can choose any k -type $P_{U,X,S}^{(b)}$ whose \mathcal{S} -marginal $P_S^{(b)}$ is $P_{\mathbf{s}^{(b)}}$, we can choose $P_{U,X,S}^{(b)} = P_{\mathbf{s}^{(b)}} \times P_{U,X|S}^{(b)}$, where $P_{U,X|S}^{(b)}$ is the conditional k -type that—among all conditional k -types—maximizes

$$\min_{\substack{P_{Y|U,X,S}: \\ P_{Y|U=u,X,S} \in \mathcal{P}(W), \forall u \in \mathcal{U}}} I(U;Y) - I(U;S), \quad (103)$$

where the mutual informations are computed w.r.t. the joint PMF $P_{\mathbf{s}^{(b)}} \times P_{U,X|S}^{(b)} \times P_{Y|U,X,S}$. Every conditional PMF can be approximated in the total variation distance by a conditional k -type when k is sufficiently large; and, because entropy and mutual information are continuous in this distance [8, Lemma 2.7], it follows that—for the above choice of the conditional k -type and some $\gamma_k = \gamma_k(|\mathcal{U}|, |\mathcal{X}|, |\mathcal{S}|, |\mathcal{Y}|)$, which converges to zero as k tends to infinity—(99) implies that when $|\mathcal{M}_{b-1}| \geq 2^{k \log |\mathcal{U}|}$

$$|\mathcal{M}_b| \leq \left(\max_{P_S} \min_{P_{U,X|S}} \max_{\substack{P_{Y|U,X,S}: \\ P_{Y|U=u,X,S} \in \mathcal{P}(W), \forall u \in \mathcal{U}}} 2^{-k(I(U;Y) - I(U;S) - \epsilon - \gamma_k)} \right) |\mathcal{M}_{b-1}|, \quad (104)$$

where the mutual informations are computed w.r.t. the joint PMF $P_S \times P_{U,X|S} \times P_{Y|U,X,S}$. Because our scheme works for any $\epsilon > 0$, it follows that for every $\epsilon > 0$ and positive integer $k \geq \eta_0(|\mathcal{X}|, |\mathcal{S}|, \epsilon)$ each of Blocks 1 through B is guaranteed to reduce the survivor set by a factor of at least

$$\min_{P_S} \max_{P_{U,X|S}} \min_{\substack{P_{Y|U,X,S}: \\ P_{Y|U=u,X,S} \in \mathcal{P}(W), \forall u \in \mathcal{U}}} 2^{k(I(U;Y) - I(U;S) - \delta(\epsilon, k))}, \quad (105)$$

until $|\mathcal{M}_B|$ is smaller than $2^{k \log |\mathcal{U}|}$. Here the mutual informations are computed w.r.t. the joint PMF $P_S \times P_{U,X|S} \times P_{Y|U,X,S}$, and

$$\delta(\epsilon, k) = \epsilon + \gamma_k \quad (106)$$

and hence converges to zero as ϵ tends to zero and k to infinity.

Since $C_{f,0}$ is positive, so is the RHS of (11); and, because $\delta(\epsilon, k)$ converges to zero as $\epsilon \downarrow 0$ and $k \rightarrow \infty$, it follows that we can choose ϵ sufficiently small and B and k sufficiently large so that

$$k \geq \eta_0(|\mathcal{X}|, |\mathcal{S}|, \epsilon) \quad (107a)$$

and

$$\left(\max_{P_S} \min_{P_{U,X|S}} \max_{\substack{P_{Y|U,X,S}: \\ P_{Y|U=u,X,S} \in \mathcal{P}(W), \forall u \in \mathcal{U}}} 2^{-Bk(I(U;Y) - I(U;S) - \delta(\epsilon, k))} \right) |\mathcal{M}| \leq 2^{k \log |\mathcal{U}|} \quad (107b)$$

This guarantees that

$$|\mathcal{M}_B| \leq 2^{k \log |\mathcal{U}|}, \quad (108)$$

because each block reduces the survivor set by the factor in (105) until $|\mathcal{M}_B|$ is smaller than $2^{k \log |\mathcal{U}|}$.

We now deal with Block $B+1$. In Block $(B+1)$ we resolve the post-Block- B survivor-set \mathcal{M}_B , and we transmit the empirical types $P_{\mathbf{s}^{(1)}}, \dots, P_{\mathbf{s}^{(B)}}$ of the state sequences pertaining to Blocks 1 through B . By (108) the resolution of \mathcal{M}_B requires at most $k \log |\mathcal{U}|$ bits. And since the empirical type of each $\mathbf{s}^{(b)}$ can take on at most $(1+k)^{|\mathcal{S}|}$ values, we need at most $B \log(1+k) |\mathcal{S}|$ bits to describe $P_{\mathbf{s}^{(1)}}, \dots, P_{\mathbf{s}^{(B)}}$. In the last block we thus need to transmit at most

$$\lceil k \log |\mathcal{U}| + B \log(1+k) |\mathcal{S}| \rceil \quad (109)$$

bits error-free. Remark 3.1 and the assumption that $C_{f,0}$ is positive guarantee that this can be achieved by choosing the length of the last block to be

$$\lceil k \log |\mathcal{U}| + B \log(1+k) |\mathcal{S}| \rceil n_{\text{bit}}, \quad (110)$$

where $n_{\text{bit}} = n_{\text{bit}}(|\mathcal{S}|, |\mathcal{Y}|)$.

We are now ready to join the dots and conclude that the coding scheme asymptotically achieves any rate smaller than the RHS of (11). More precisely, we will show that, for every rate R smaller than the RHS of (11) and every sufficiently-large blocklength n , our coding scheme can convey nR bits error-free in n channel uses. It follows from (107) and (110) that if the positive integers n , B , k and $\epsilon > 0$ are such that

$$k \geq \eta_0(|\mathcal{X}|, |\mathcal{S}|, \epsilon) \quad (111a)$$

and

$$nR \leq Bk \left(\min_{P_S} \max_{P_{U,X|S}} \min_{\substack{P_{Y|U,X,S}: \\ P_{Y|U=u,X,S} \in \mathcal{P}(W), \forall u \in \mathcal{U}}} I(U; Y) - I(S; Y) - \delta(\epsilon, k) \right), \quad (111b)$$

then our coding scheme can convey nR bits error-free in

$$Bk + \lceil k \log |\mathcal{U}| + B \log(1+k) |\mathcal{S}| \rceil n_{\text{bit}} \quad (112)$$

channel uses. It thus remains to exhibit positive integers B , k and some $\epsilon > 0$ such that, for every sufficiently-large blocklength n , (111) holds and

$$Bk + \lceil k \log |\mathcal{U}| + B \log(1+k) |\mathcal{S}| \rceil n_{\text{bit}} \leq n. \quad (113)$$

As we argue next, when n is sufficiently large we can choose

$$B = \lfloor \sqrt{n} \rfloor - \lceil \log |\mathcal{U}| + \log(1 + \sqrt{n}) |\mathcal{S}| \rceil n_{\text{bit}}, \quad (114a)$$

$$k = \lfloor \sqrt{n} \rfloor, \quad (114b)$$

and we can choose any $\epsilon > 0$ for which

$$R + \epsilon < \min_{P_S} \max_{P_{U,X|S}} \min_{\substack{P_{Y|U,X,S}: \\ P_{Y|U=u,X,S} \in \mathcal{P}(W), \forall u \in \mathcal{U}}} I(U; Y) - I(S; Y). \quad (115)$$

Note that, whenever n is sufficiently large, B is positive and (113) is satisfied. To see that also (111) holds whenever n is sufficiently large, we first observe from (114b) that k tends to infinity as n tends to infinity. This implies that (111a) holds whenever n is sufficiently large, and that $\delta(\epsilon, k)$ (which is defined in (106), where $\gamma_k = \gamma_k(|\mathcal{U}|, |\mathcal{X}|, |\mathcal{S}|, |\mathcal{Y}|)$ converges to zero as k tends to infinity) converges to ϵ as n tends to infinity. We next observe that (114) implies that Bk/n converges to one as n tends to infinity. This, combined with the fact that $\delta(\epsilon, k)$ converges to ϵ as n tends to infinity and with (115), implies that (111b) holds whenever n is sufficiently large. \square

We next prove the converse part of Theorem 2.4.

Converse Part. Fix a finite set \mathcal{M} , a blocklength n , and an (n, \mathcal{M}) zero-error code with n encoding mappings

$$f_i: \mathcal{M} \times \mathcal{S}^n \times \mathcal{Y}^{i-1} \rightarrow \mathcal{X}, \quad i \in [1 : n] \quad (116)$$

and $|\mathcal{M}|$ disjoint decoding sets $\mathcal{D}_m \subseteq \mathcal{Y}^n$, $m \in \mathcal{M}$. We will show that, for some chance variable U of finite support \mathcal{U} , the rate $\frac{1}{n} \log |\mathcal{M}|$ of the code is upper-bounded by the RHS of (11).

Draw M uniformly over \mathcal{M} , and denote its distribution P_M . Since the code is a zero-error code,

$$\mathbb{P}[Y^n \in \mathcal{D}_M] = 1, \quad (117)$$

where \mathbb{P} is the distribution of (M, S^n, X^n, Y^n) induced by P_M , the state distribution Q , the encoding mappings (116), and the channel law $W(y|x, s)$, so for every $(m, \mathbf{s}, \mathbf{x}, \mathbf{y}) \in \mathcal{M} \times \mathcal{S}^n \times \mathcal{X}^n \times \mathcal{Y}^n$

$$\begin{aligned} \mathbb{P}[(M, S^n, X^n, Y^n) = (m, \mathbf{s}, \mathbf{x}, \mathbf{y})] \\ = P_M(m) Q^n(\mathbf{s}) \prod_{i=1}^n \left(P_{X_i|M, S^n, Y^{i-1}}(x_i|m, \mathbf{s}, y^{i-1}) W(y_i|x_i, s_i) \right), \end{aligned} \quad (118)$$

where

$$P_{X_i|M, S^n, Y^{i-1}}(x_i|m, \mathbf{s}, y^{i-1}) = \begin{cases} 1 & \text{if } x_i = f_i(m, \mathbf{s}, y^{i-1}), \\ 0 & \text{otherwise.} \end{cases} \quad (119)$$

Fix any PMF \tilde{P}_S on \mathcal{S} and any collection of n conditional PMFs $\{\tilde{P}_{Y_i|M, Y^{i-1}, S_{i+1}^n, X_i, S_i}\}_{i \in [1:n]}$ that satisfy

$$\begin{aligned} \tilde{P}_{Y_i|M, Y^{i-1}, S_{i+1}^n, X_i, S_i}(\cdot|m, y^{i-1}, s_{i+1}^n, x_i, s_i) &\ll W(\cdot|x_i, s_i), \\ \forall (m, y^{i-1}, s_{i+1}^n, x_i, s_i) &\in \mathcal{M} \times \mathcal{Y}^{i-1} \times \mathcal{S}^{n-i} \times \mathcal{X} \times \mathcal{S}. \end{aligned} \quad (120)$$

These PMFs induce the PMF on $\mathcal{M} \times \mathcal{S}^n \times \mathcal{X}^n \times \mathcal{Y}^n$

$$\tilde{P}_{M, S^n, X^n, Y^n} = P_M \times \tilde{P}_S^n \times \prod_{i=1}^n (P_{X_i|M, S^n, Y^{i-1}} \times \tilde{P}_{Y_i|M, Y^{i-1}, S_{i+1}^n, X_i, S_i}). \quad (121)$$

It follows from (1) and (120) that $\tilde{P}_{M, S^n, X^n, Y^n} \ll \mathbb{P}$ and consequently that (117) implies

$$\tilde{P}_{M, S^n, X^n, Y^n}[Y^n \in \mathcal{D}_M] = 1. \quad (122)$$

We upper-bound $\frac{1}{n} \log |\mathcal{M}|$ by carrying out the following calculation as in [12, Section 7.6] but under $\tilde{P}_{M,S^n,X^n,Y^n}$ of (121):

$$\begin{aligned} & \frac{1}{n} \log |\mathcal{M}| \\ & \stackrel{(a)}{=} \frac{1}{n} H(M) \end{aligned} \tag{123}$$

$$\stackrel{(b)}{=} \frac{1}{n} \left[I(M; Y^n) - I(M; S^n) \right] \tag{124}$$

$$\stackrel{(c)}{=} \frac{1}{n} \sum_{i=1}^n \left[I(M; Y_i | Y^{i-1}) - I(M; S_i | S_{i+1}^n) \right] \tag{125}$$

$$\begin{aligned} & \stackrel{(d)}{=} \frac{1}{n} \sum_{i=1}^n \left[I(M, Y^{i-1}, S_{i+1}^n; Y_i) - I(Y^{i-1}; Y_i) - I(S_{i+1}^n; Y_i | M, Y^{i-1}) \right. \\ & \quad \left. - I(M, Y^{i-1}, S_{i+1}^n; S_i) + I(S_{i+1}^n; S_i) + I(Y^{i-1}; S_i | M, S_{i+1}^n) \right] \end{aligned} \tag{126}$$

$$\stackrel{(e)}{\leq} \frac{1}{n} \sum_{i=1}^n \left[I(M, Y^{i-1}, S_{i+1}^n; Y_i) - I(M, Y^{i-1}, S_{i+1}^n; S_i) \right], \tag{127}$$

where (a) holds because under $\tilde{P}_{M,S^n,X^n,Y^n}$ M is uniform over \mathcal{M} ; (b) holds by (122) and because under $\tilde{P}_{M,S^n,X^n,Y^n}$ M is independent of S^n ; (c) and (d) follow from the chain rule; and (e) follows from Csiszár's sum-identity, the nonnegativity of mutual information, and the independence of S_i and S_{i+1}^n under $\tilde{P}_{M,S^n,X^n,Y^n}$.

For every $i \in [1 : n]$ define the chance variable

$$U_i = (M, Y^{i-1}, S_{i+1}^n). \tag{128}$$

From (127) it then follows that every choice of \tilde{P}_S and $\{\tilde{P}_{Y_i|M,Y^{i-1},S_{i+1}^n,X_i,S_i}\}_{i \in [1:n]}$ satisfying (120) gives rise to an upper bound

$$\frac{1}{n} \log |\mathcal{M}| \leq \frac{1}{n} \sum_{i=1}^n \left[I(U_i; Y_i) - I(U_i; S_i) \right], \tag{129}$$

where the mutual informations in the i -th summand are computed w.r.t. the joint PMF $\tilde{P}_{U_i,X_i,S_i,Y_i}$ induced by $\tilde{P}_{M,S^n,X^n,Y^n}$.

We will conclude the proof by exhibiting a PMF \tilde{P}_S and a collection of conditional PMFs $\{\tilde{P}_{Y_i|M,Y^{i-1},S_{i+1}^n,X_i,S_i}\}_{i \in [1:n]}$ satisfying (120) for which each summand on the RHS of (129) is upper-bounded by the RHS of (11).

We begin with the choice of $\{\tilde{P}_{Y_i|M,Y^{i-1},S_{i+1}^n,X_i,S_i}\}_{i \in [1:n]}$. To this end note from (128) the one-to-one correspondence between $\tilde{P}_{Y_i|M,Y^{i-1},S_{i+1}^n,X_i,S_i}$ and $\tilde{P}_{Y_i|U_i,X_i,S_i}$:

$$\begin{aligned} & \tilde{P}_{Y_i|M,Y^{i-1},S_{i+1}^n,X_i,S_i}(y_i|m, y^{i-1}, s_{i+1}^n, x_i, s_i) = \tilde{P}_{Y_i|U_i,X_i,S_i}(y_i|(m, y^{i-1}, s_{i+1}^n), x_i, s_i), \\ & \forall (m, y^{i-1}, s_{i+1}^n, x_i, s_i) \in \mathcal{M} \times \mathcal{Y}^{i-1} \times \mathcal{S}^{n-i} \times \mathcal{X} \times \mathcal{S}. \end{aligned} \tag{130}$$

This implies that choosing a conditional PMF $\tilde{P}_{Y_i|M, Y^{i-1}, S_{i+1}^n, X_i, S_i}$ that satisfies (120) is tantamount to choosing a conditional PMF $\tilde{P}_{Y_i|U_i, X_i, S_i}$ that satisfies

$$\tilde{P}_{Y_i|U_i=u_i, X_i, S_i} \in \mathcal{P}(W), \quad \forall u_i \in \mathcal{U}_i, \quad (131)$$

and consequently choosing a collection of conditional PMFs $\{\tilde{P}_{Y_i|M, Y^{i-1}, S_{i+1}^n, X_i, S_i}\}_{i \in [1:n]}$ that satisfy (120) is tantamount to choosing a collection of conditional PMFs $\{\tilde{P}_{Y_i|U_i, X_i, S_i}\}_{i \in [1:n]}$ that satisfy (131). We shall choose the latter collection, and we shall do so as follows.

We first choose $\tilde{P}_{Y_i|U_i, X_i, S_i}$ for $i = 1$, and we then repeatedly increment i by one until it reaches n . Key to our choice is the observation, which will be justified shortly, that $\tilde{P}_{U_i, X_i, S_i}$ is determined by \tilde{P}_S and $\{\tilde{P}_{Y_j|U_j, X_j, S_j}\}_{j \in [1:i-1]}$. Our choice of $\tilde{P}_{Y_i|U_i, X_i, S_i}$ can thus depend not only on our choice of \tilde{P}_S and our previous choices of $\{\tilde{P}_{Y_j|U_j, X_j, S_j}\}_{j \in [1:i-1]}$ but also on $\tilde{P}_{U_i, X_i, S_i}$. This will allow us to choose $\tilde{P}_{Y_i|U_i, X_i, S_i}$ as one that—among all conditional PMFs satisfying (131)—minimizes

$$I(U_i; Y_i) - I(U_i; S_i), \quad (132)$$

where the mutual informations are computed w.r.t. the joint PMF $\tilde{P}_{U_i, X_i, S_i} \times \tilde{P}_{Y_i|U_i, X_i, S_i}$. Since (121) implies that

$$\tilde{P}_{S_i} = \tilde{P}_S, \quad i \in [1 : n], \quad (133)$$

we will then find that, for our choice of $\{\tilde{P}_{Y_i|U_i, X_i, S_i}\}_{i \in [1:n]}$,

$$\begin{aligned} & I(U_i; Y_i) - I(U_i; S_i) \\ & \leq \max_{\tilde{P}_{U_i, X_i|S_i}} \min_{\substack{\tilde{P}_{Y_i|U_i, X_i, S_i} : \\ \tilde{P}_{Y_i|U_i=u_i, X_i, S_i} \in \mathcal{P}(W), \forall u_i \in \mathcal{U}_i}} I(U_i; Y_i) - I(U_i; S_i), \quad i \in [1 : n], \end{aligned} \quad (134)$$

where the mutual informations are computed w.r.t. the joint PMF $\tilde{P}_{S_i} \times \tilde{P}_{U_i, X_i|S_i} \times \tilde{P}_{Y_i|U_i, X_i, S_i}$. The chosen conditional PMFs $\{\tilde{P}_{Y_i|U_i, X_i, S_i}\}_{i \in [1:n]}$ satisfy (131), and hence (129) and (134) will imply that

$$\begin{aligned} & \frac{1}{n} \log |\mathcal{M}| \\ & \leq \frac{1}{n} \sum_{i=1}^n \max_{\tilde{P}_{U_i, X_i|S_i}} \min_{\substack{\tilde{P}_{Y_i|U_i, X_i, S_i} : \\ \tilde{P}_{Y_i|U_i=u_i, X_i, S_i} \in \mathcal{P}(W), \forall u_i \in \mathcal{U}_i}} I(U_i; Y_i) - I(U_i; S_i), \end{aligned} \quad (135)$$

where the mutual informations in the i -th summand are computed w.r.t. $\tilde{P}_{S_i} \times \tilde{P}_{U_i, X_i|S_i} \times \tilde{P}_{Y_i|U_i, X_i, S_i}$.

We now prove that indeed $\tilde{P}_{U_i, X_i, S_i}$ is determined by \tilde{P}_S and $\{\tilde{P}_{Y_j|U_j, X_j, S_j}\}_{j \in [1:i-1]}$. In fact, we will show that the latter two determine $\tilde{P}_{M, S^n, X^i, Y^{i-1}}$. The latter determines $\tilde{P}_{U_i, X_i, S_i}$, because, by (128), the tuple (U_i, X_i, S_i) is determined by (M, S^n, X^i, Y^{i-1}) and consequently its PMF $\tilde{P}_{U_i, X_i, S_i}$ is determined by $\tilde{P}_{M, S^n, X^i, Y^{i-1}}$.

We use mathematical induction, but first we note that the PMF $\tilde{P}_{M, S^n, X^n, Y^n}$ is constructed inductively: by (121)

$$\tilde{P}_{M, S^n, X_1} = P_M \times \tilde{P}_S^n \times P_{X_1|M, S^n} \quad (136)$$

and, for every $\ell \in [2 : n]$, $\tilde{P}_{M, S^n, X^\ell, Y^{\ell-1}}$ is constructed from $\tilde{P}_{M, S^n, X^{\ell-1}, Y^{\ell-2}}$ by

$$\tilde{P}_{M, S^n, X^\ell, Y^{\ell-1}} = \tilde{P}_{M, S^n, X^{\ell-1}, Y^{\ell-2}} \times \tilde{P}_{Y_{\ell-1}|M, Y^{\ell-2}, S_\ell^n, X_{\ell-1}, S_{\ell-1}} \times P_{X_\ell|M, S^n, Y^{\ell-1}}. \quad (137)$$

In describing the proof we shall make the dependence on P_M , our choice of \tilde{P}_S , and $\{P_{X_j|M, S^n, Y^{j-1}}\}_{j \in [1:n]}$, whose components are determined by the encoding mappings (116) via (119), implicit.

1. Basis $\ell = 1$: It follows from (136) that \tilde{P}_{M, S^n, X_1} is determined.
2. Inductive Step: Fix $\ell \in [2 : i]$, and suppose that $\tilde{P}_{M, S^n, X^{\ell-1}, Y^{\ell-2}}$ is determined by $\{\tilde{P}_{Y_j|U_j, X_j, S_j}\}_{j \in [1:\ell-2]}$. Since $\tilde{P}_{Y_{\ell-1}|M, Y^{\ell-2}, S_\ell^n, X_{\ell-1}, S_{\ell-1}}$ is by (128) in a one-to-one correspondence with $\tilde{P}_{Y_{\ell-1}|U_{\ell-1}, X_{\ell-1}, S_{\ell-1}}$, this implies that $\tilde{P}_{M, S^n, X^{\ell-1}, Y^{\ell-2}}$ and $\tilde{P}_{Y_{\ell-1}|M, Y^{\ell-2}, S_\ell^n, X_{\ell-1}, S_{\ell-1}}$ are determined by $\{\tilde{P}_{Y_j|U_j, X_j, S_j}\}_{j \in [1:\ell-1]}$. Consequently, it follows from (137) that $\tilde{P}_{M, S^n, X^\ell, Y^{\ell-1}}$ is determined by $\{\tilde{P}_{Y_j|U_j, X_j, S_j}\}_{j \in [1:\ell-1]}$.

This proves that, for every $i \in [1 : n]$, $\tilde{P}_{M, S^n, X^i, Y^{i-1}}$ and consequently also $\tilde{P}_{U_i, X_i, S_i}$ are determined by \tilde{P}_S and $\{\tilde{P}_{Y_j|U_j, X_j, S_j}\}_{j \in [1:i-1]}$, and hence (135) holds.

Having established (135), we are now ready to conclude the proof. By the definition of U_i (128) the cardinality of the support \mathcal{U}_i of U_i satisfies

$$|\mathcal{U}_i| \leq |\mathcal{M}| \max\{|\mathcal{Y}|, |\mathcal{S}|\}^n, \quad i \in [1 : n]. \quad (138)$$

Consequently, (133) and (135) imply that

$$\frac{1}{n} \log |\mathcal{M}| \leq \max_{\tilde{P}_{U, X|S}} \min_{\substack{\tilde{P}_{Y|U, X, S}: \\ \tilde{P}_{Y|U=u, X, S} \in \mathcal{P}(W), \forall u \in \mathcal{U}}} I(U; Y) - I(U; S), \quad (139)$$

where U is an auxiliary chance variable taking values in a finite set \mathcal{U} , and the mutual informations are computed w.r.t. the joint PMF $\tilde{P}_S \times \tilde{P}_{U, X|S} \times \tilde{P}_{Y|U, X, S}$. Since we can choose any PMF \tilde{P}_S on \mathcal{S} , we can choose one that—among all PMFs on \mathcal{S} —yields the tightest bound, i.e., minimizes

$$\max_{\tilde{P}_{U, X|S}} \min_{\substack{\tilde{P}_{Y|U, X, S}: \\ \tilde{P}_{Y|U=u, X, S} \in \mathcal{P}(W), \forall u \in \mathcal{U}}} I(U; Y) - I(U; S), \quad (140)$$

where U is an auxiliary chance variable taking values in a finite set \mathcal{U} , and the mutual informations are computed w.r.t. the joint PMF $\tilde{P}_S \times \tilde{P}_{U,X|S} \times \tilde{P}_{Y|U,X,S}$. For this choice of \tilde{P}_S (139) implies that

$$\frac{1}{n} \log |\mathcal{M}| \leq \min_{\tilde{P}_S} \max_{\tilde{P}_{U,X|S}} \min_{\substack{\tilde{P}_{Y|U,X,S}: \\ \tilde{P}_{Y|U=u,X,S} \in \mathcal{P}(W), \forall u \in \mathcal{U}}} I(U; Y) - I(U; S), \quad (141)$$

where U is an auxiliary chance variable taking values in a finite set \mathcal{U} , and the mutual informations are computed w.r.t. the joint PMF $\tilde{P}_S \times \tilde{P}_{U,X|S} \times \tilde{P}_{Y|U,X,S}$. \square

3.3 A Proof of Theorem 2.7

We use the following lemma to establish Theorem 2.7:

Lemma 3.3 (No Feedback). *In the absence of feedback, a sufficient condition for the zero-error capacity of the SD-DMC $W(y|x, s)$ with acausal SI to be zero is*

$$\exists s \in \mathcal{S} \quad \forall x \in \mathcal{X} \quad \exists s' \in \mathcal{S} \quad \forall x' \in \mathcal{X} \quad \exists y \in \mathcal{Y} \text{ s.t. } W(y|x, s) W(y|x', s') > 0. \quad (142)$$

A sufficient condition for the capacity in the absence of feedback to be positive is that for some $\kappa \in [2 : |\mathcal{Y}|]$ and $\lambda \in [2 : \kappa |\mathcal{X}|]$ there exist channel inputs

$$x(s, k), \quad (s, k) \in \mathcal{S} \times [1 : \kappa]$$

and λ pairwise-disjoint nonempty subsets $\mathcal{Y}_1, \dots, \mathcal{Y}_\ell \subset \mathcal{Y}$ such that the following two conditions hold:

$$\forall (s, k) \in \mathcal{S} \times [1 : \kappa] \quad \exists \ell \in [1 : \lambda] \text{ s.t. } W(\mathcal{Y}_\ell | x(s, k), s) = 1 \quad (143a)$$

and

$$\forall \ell \in [1 : \lambda] \quad \exists k' \in [1 : \kappa] \text{ s.t. } \left(W(\mathcal{Y}_\ell | x(s', k'), s') = 0, \forall s' \in \mathcal{S} \right). \quad (143b)$$

Proof. We first prove that if (142) holds, then without feedback it is impossible to transmit a single bit error-free. Let the bit take values in the set $\mathcal{M} = \{0, 1\}$, and fix a blocklength n , an encoding mapping $f: \mathcal{M} \times \mathcal{S}^n \rightarrow \mathcal{X}^n$, and two disjoint decoding sets $\mathcal{D}_m \subseteq \mathcal{Y}^n$, $m \in \mathcal{M}$. By (142) there exists some state $s^* \in \mathcal{S}$ for which

$$\forall x \in \mathcal{X} \quad \exists s' \in \mathcal{S} \quad \forall x' \in \mathcal{X} \quad \exists y \in \mathcal{Y} \text{ s.t. } W(y|x, s^*) W(y|x', s') > 0. \quad (144)$$

Let $\mathbf{s}^* \in \mathcal{S}^n$ be the all- s^* state-sequence, so $s_i^* = s^*$, $i \in [1 : n]$, and let $\mathbf{x} = f(0, \mathbf{s}^*)$. Choosing x in (144) to be the i -th component x_i of $f(0, \mathbf{s}^*)$, it follows from (144) that for every $i \in [1 : n]$ there exists some $s' \in \mathcal{S}$, say $s'(i)$, for which

$$\forall x' \in \mathcal{X} \quad \exists y \in \mathcal{Y} \text{ s.t. } W(y|x_i, s_i^*) W(y|x', s'(i)) > 0. \quad (145)$$

Let $\mathbf{s}' \in \mathcal{S}^n$ be the state sequence whose i -th component s'_i is $s'(i)$, $i \in [1 : n]$, and let $\mathbf{x}' = f(1, \mathbf{s}')$. By (145)

$$\exists \mathbf{y} \in \mathcal{Y}^n \text{ s.t. } \prod_{i=1}^n \left(W(y_i | x_i, s_i^*) W(y_i | x'_i, s'_i) \right) > 0. \quad (146)$$

This makes it impossible for the decoder to determine with certainty whether the transmitted bit is 0 or 1 even if it is told that the state sequence is \mathbf{s}^* or \mathbf{s}' . This concludes the proof of the first part of the lemma.

It remains to prove that if for some $\kappa \in [2 : |\mathcal{Y}|]$ and $\lambda \in [2 : \kappa |\mathcal{X}|]$ there exist channel inputs $\{x(s, k)\}_{(s,k) \in \mathcal{S} \times [1:\kappa]}$ and pairwise-disjoint output-sets $\{\mathcal{Y}_\ell\}_{\ell \in [1:\lambda]}$ for which (143) holds, then the no-feedback zero-error capacity of the SD-DMC $W(y|x, s)$ with acausal SI is positive. The proof is similar to that of Remark 3.1. To make up for the missing feedback, we shall choose the inputs so that the encoder—while incognizant of Y_i —will know which of the subsets $\{\mathcal{Y}_\ell\}_{\ell \in [1:\lambda]}$ contains Y_i . The decoder will, of course, know that too.

If there is only one state s^* , i.e., $\mathcal{S} = \{s^*\}$, then upon defining $x \triangleq x(s^*, 1)$ we obtain from (143a) the existence of some $\ell \in [1 : \lambda]$ for which

$$W(\mathcal{Y}_\ell | x, s^*) = 1. \quad (147a)$$

It then follows from (143b) that there exists some $k' \in [1 : \kappa]$ with corresponding $x' = x(s^*, k')$ for which

$$W(\mathcal{Y}_\ell | x', s^*) = 0. \quad (147b)$$

From (147) we obtain that

$$W(y | x, s^*) W(y | x', s^*) = 0, \quad \forall y \in \mathcal{Y}, \quad (148)$$

and by sending x or x' we can transmit a bit error-free. We hence consider now $|\mathcal{S}| \geq 2$.

To transmit a single bit $m \in \{0, 1\}$, we use two phases of n_1 and n_2 channel uses, where

$$n_{\text{bit}} = n_1 + n_2. \quad (149)$$

The goal of Phase 1 is to produce a random subset $\mathcal{S}_{n_1} \subseteq \mathcal{S}^{n_1}$ with the following three properties: 1) both encoder and decoder know \mathcal{S}_{n_1} before Phase 2 begins; 2) with probability one \mathcal{S}_{n_1} contains the Phase-2 state-sequence $\mathcal{S}_{n_1+1}^{n_1+n_2}$; and 3) the cardinality of \mathcal{S}_{n_1} is upper-bounded by

$$|\mathcal{S}_{n_1}| \leq \left(\frac{\kappa - 1}{\kappa} \right)^{n_1} |\mathcal{S}|^{n_2} + \kappa. \quad (150)$$

To that end we partition the set $\mathcal{S}_0 = \mathcal{S}^{n_2}$ into κ different subsets whose size is between $\lfloor |\mathcal{S}_0|/\kappa \rfloor$ and $\lceil |\mathcal{S}_0|/\kappa \rceil$. We index the κ subsets by the set $[1 : \kappa]$ and reveal the result to the encoder and decoder. If, thanks to its acausal SI, the encoder knows that the Time-1 state S_1 is s and that $S_{n_1+1}^{n_1+n_2}$ is in the subset of \mathcal{S}_0 indexed by k , then at Time 1 it transmits $x(s, k)$. By (143a) there exists some $\ell^* \in [1 : \lambda]$ with corresponding subset \mathcal{Y}_{ℓ^*} such that, with probability one, Y_1 is in \mathcal{Y}_{ℓ^*} . And since the subsets $\{\mathcal{Y}_{\ell}\}_{\ell \in [1:\lambda]}$ are pairwise disjoint, the probability of Y_1 being in another subset is zero. The decoder can thus compute ℓ^* from Y_1 by checking which subset contains Y_1 . The encoder knows ℓ^* , because it knows the pair (s, k) . Based on \mathcal{Y}_{ℓ^*} the encoder and decoder can determine all $k' \in [1 : \kappa]$ for which

$$W(\mathcal{Y}_{\ell^*} | x(s', k'), s') = 0, \forall s' \in \mathcal{S}. \quad (151)$$

(By (143b) at least one such k' exists.) Because $Y_1 \in \mathcal{Y}_{\ell^*}$ and by (151), the Phase-2 state-sequence cannot be contained in a subset of \mathcal{S}_0 indexed by such a k' , and hence it is in the \mathcal{S}_0 -complement of these subsets, which we denote \mathcal{S}_1 . Note that: 1) both encoder and decoder know \mathcal{S}_1 after Channel-Use 1; 2) \mathcal{S}_1 contains $S_{n_1+1}^{n_1+n_2}$; and 3) the cardinality of \mathcal{S}_1 is upper-bounded by

$$|\mathcal{S}_1| \leq |\mathcal{S}_0| - \left\lfloor \frac{|\mathcal{S}_0|}{\kappa} \right\rfloor \leq \frac{\kappa - 1}{\kappa} |\mathcal{S}_0| + 1. \quad (152)$$

Phase 1 continues in the same fashion, and hence we obtain that, for every $i \in [1 : n_1]$, the first i channel uses produce a random subset \mathcal{S}_i of \mathcal{S}^{n_2} satisfying that: 1) both encoder and decoder know \mathcal{S}_i after Channel-Use i ; 2) \mathcal{S}_i contains $S_{n_1+1}^{n_1+n_2}$; and 3) the cardinality of \mathcal{S}_i is upper-bounded by

$$|\mathcal{S}_i| \leq |\mathcal{S}_{i-1}| - \left\lfloor \frac{|\mathcal{S}_{i-1}|}{\kappa} \right\rfloor \leq \frac{\kappa - 1}{\kappa} |\mathcal{S}_{i-1}| + 1. \quad (153)$$

As in the proof of Remark 3.1, this implies that Phase 1 produces a random subset \mathcal{S}_{n_1} of \mathcal{S}^{n_2} with the desired three properties.

Phase 2 in the proof of Remark 3.1 does not use the feedback link, and hence we can use it also in the current setting without feedback. Consequently, we can argue essentially as in the proof of Remark 3.1 but with (77) replaced by

$$n_1 = \left\lceil \frac{2\kappa \log |\mathcal{S}| - \log \kappa}{\log \kappa - \log(\kappa - 1)} \right\rceil, \quad (154a)$$

$$n_2 = 2\kappa, \quad (154b)$$

$$n_{\text{bit}} = \left\lceil \frac{2\kappa \log |\mathcal{S}| - \log \kappa}{\log \kappa - \log(\kappa - 1)} \right\rceil + 2\kappa \quad (154c)$$

that n_{bit} channel uses suffice for the error-free transmission of a single bit. This concludes the proof, because κ is at most $|\mathcal{Y}|$ and hence it follows from (154) that n_{bit} satisfies the upper bound

$$n_{\text{bit}} \leq \left\lceil \frac{2|\mathcal{Y}| \log |\mathcal{S}| - \log |\mathcal{Y}|}{\log |\mathcal{Y}| - \log(|\mathcal{Y}| - 1)} \right\rceil + 2|\mathcal{Y}|. \quad (155)$$

□

Theorem 2.7 follows from Theorem 2.3, Lemma 3.3, and the following example:

Example 3.4. Suppose $\mathcal{X} = \{0, 1\}$ and $\mathcal{S} = \mathcal{Y} = \{1, 2, 3, 4, 5\}$. For every $x \in \mathcal{X}$ and $s \in \mathcal{S}$ define $\mathcal{Y}_{x,s}$ according to Table 3, and let $W(y|x, s)$ be such that

$$\{y \in \mathcal{Y} : W(y|x, s) > 0\} = \mathcal{Y}_{x,s}, \quad \forall (x, s) \in \mathcal{X} \times \mathcal{S}. \quad (156)$$

Then, the SD-DMC $W(y|x, s)$ satisfies both (10) and (142).

$\mathcal{Y}_{x,s}$		s				
		1	2	3	4	5
x	0	{ 2,3 }	{ 1,5 }	{ 1,2 }	{ 2,3 }	{ 1,2 }
	1	{ 4,5 }	{ 3,4 }	{ 4,5 }	{ 1,5 }	{ 3,4 }

Table 3: Nonzero transitions of the SD-DMC in Example 3.4.

Remark 3.5. Lemma 3.3 does not fully characterize the SD-DMCs whose capacity is positive in the absence of feedback. For example the SD-DMC of Example 3.4 but with state alphabet $\mathcal{S} = \{1, 2, 4\}$ satisfies neither the conditions of the lemma. However, when $W(y|x, s)$ is $\{0, 1\}$ -valued (cf. Example 2.9), Lemma 3.3 implies that the capacity is positive iff

$$|\{y \in \mathcal{Y} : \exists x \in \mathcal{X} \text{ s.t. } W(y|x, s) > 0\}| \geq 2, \quad \forall s \in \mathcal{S}. \quad (157)$$

(To see this, choose the sets $\{\mathcal{Y}_\ell\}$ in Lemma 3.3 to be the singletons containing the outputs $y \in \mathcal{Y}$ for which $W(y|x, s) > 0$ holds for some $(x, s) \in \mathcal{X} \times \mathcal{S}$.)

4 Summary

We now know the zero-error feedback capacity of the state-dependent channel in all three cases: when the state is revealed to the encoder strictly-causally, causally, or acausally. In each case the capacity result comprises two parts: a characterization of the channels for which the capacity is positive, and a formula for the capacity when it is.

- Revealing the state to the encoder strictly-causally does not increase capacity (Remark 2.17), and the problem reduces to the state-less channel, which was solved by Shannon [1], with Ahlswede [5] later providing an alternative form and an alternative blocks-based coding scheme.
- When the state is revealed to the encoder causally, the SI is utilized optimally by using Shannon strategies, and the zero-error feedback capacity is thus that of the state-less channel into which the state-dependent channel is transformed when the encoder uses Shannon strategies (Theorems 2.10 and 2.11).
- For the case where the state is revealed to the encoder acausally, our positivity characterization (Theorem 2.3) is reminiscent of Shannon's, and our formula (Theorem 2.4) is reminiscent of Ahlswede's.

The acausal case exhibits phenomena that are not observed in the strictly-causal and causal cases: The zero-error feedback capacity can be positive even if in the absence of feedback the zero-error capacity is zero (Theorem 2.7), and the error-free transmission of a single bit may require more than one channel use (Corollary 2.8).

Our coding scheme for the acausal case builds on Ahlswede's blocks-based scheme [5] and to a lesser degree on Shannon's sequential approach [1]. In contrast to Shannon's sequential scheme, in Ahlswede's scheme the encoder codes over blocks, and it can therefore take advantage of the acausal SI in a more natural way. Ahlswede's scheme also seems to be more natural in the state-less case in the presence of input constraints: his expression remains valid provided we replace the maximization over the input distribution with a constrained maximization (Corollary 2.21). This is not the case for Shannon's expression (Remark 2.22).

For the acausal case we also established the zero-error feedback capacity for a scenario where—in addition to the message—also the state sequence must be recovered (Theorem 2.19); for a scenario with an average-cost constraint on the channel inputs (Theorem 2.20); and for a scenario with an average-cost constraint on prespecified l -blocks of consecutive channel states (Theorem 2.26).

A recurring theme in our coding schemes is that, as of the beginning of the transmission, the encoder attempts to convey not only the message but also the state sequence governing the last block, a state sequence of which it is cognizant because the entire state sequence is revealed to it acausally. Once the ambiguity about the last-block's state sequence and the message has been sufficiently reduced, the last block is used to resolve it, or rather to decode the message.

Another recurring theme in our coding schemes is that—to reduce the decoder's ambiguity about the message and the last-block's state sequence—each block uses pairwise

disjoint bins that “completely cover” the set of possible state sequences in the sense that all the state sequences pertaining to the block can be accommodated.

A recurring theme in the converse parts is to select the “worst possible” joint distribution of the message, state sequence, input sequence, and output sequence. By “possible” we mean here that the distribution is compatible with the encoding mappings and absolutely continuous w.r.t. the distribution that is induced by the uniform message distribution, the state distribution, the encoding mappings, and the channel law. By “worst” we mean that the distribution yields—among all “possible” distributions—the tightest bound.

A remaining open problem is to characterize the family of channels whose zero-error capacity with acausal SI is zero *in the absence of feedback*. We provided a sufficient condition (Lemma 3.3), which we then used to show that some members of this family have positive zero-error capacity in the presence of feedback (Theorem 2.7). We also showed that some channels outside this family have zero zero-error capacity when the state is revealed causally (Theorem 2.14). On such channels with acausal SI the error-free transmission of a single bit requires more than one channel use also in the absence of feedback (Corollary 2.15). (Recall that in the causal case the zero-error capacity—both in the presence and in the absence of feedback—is positive iff it is possible to transmit a single bit error-free in one channel use.) One way to characterize the family might be to upper-bound the maximal number of channel uses that could be necessary to transmit a single bit error-free.

A A Proof of Remark 2.2

Definition A.1. For any finite set \mathcal{M} and positive integer $n \in \mathbb{N}$, an (n, \mathcal{M}) zero-error feedback code with acausal SI and a stochastic encoder is defined like its deterministic counterpart (Definition 2.1) except that the encoding may depend on some chance variable Θ that is drawn from some finite set \mathcal{T} according to some PMF P_Θ .⁶ The code thus consists of a finite set \mathcal{T} , a PMF P_Θ on \mathcal{T} , n encoding mappings

$$f_i: \mathcal{M} \times \mathcal{T} \times \mathcal{S}^n \times \mathcal{Y}^{i-1} \rightarrow \mathcal{X}, \quad i \in [1 : n], \quad (158)$$

and $|\mathcal{M}|$ disjoint decoding sets

$$\mathcal{D}_m \subseteq \mathcal{Y}^n, \quad m \in \mathcal{M} \quad (159)$$

⁶The assumption that Θ takes values in a finite set is not restrictive, because the channel-input, -state, and -output alphabets are finite (see Remark A.3 at the end of this section).

such that for every $m \in \mathcal{M}$ the probability of a decoding error is zero, i.e.,

$$\mathbb{P}[Y^n \notin \mathcal{D}_m | M = m, S^n = \mathbf{s}] = 0, \quad \forall m \in \mathcal{M}, \mathbf{s} \in \mathcal{S}^n, \quad (160)$$

where

$$\begin{aligned} & \mathbb{P}[Y^n \notin \mathcal{D}_m | M = m, S^n = \mathbf{s}] \\ &= \sum_{\theta \in \mathcal{T}} P_{\Theta}(\theta) \sum_{\mathbf{y} \in \mathcal{Y}^n \setminus \mathcal{D}_m} \prod_{i=1}^n W(y_i | f_i(m, \theta, \mathbf{s}, y^{i-1}), s_i). \end{aligned} \quad (161)$$

Proof of Remark 2.2. Given an (n, \mathcal{M}) zero-error feedback code with a stochastic encoder (158) and decoding sets (159), we can construct an (n, \mathcal{M}) zero-error feedback code with a deterministic encoder (6) as follows. We fix some element θ^* of \mathcal{T} for which $P_{\Theta}(\theta^*) > 0$ and consider the n deterministic encoding mappings

$$g_i: \mathcal{M} \times \mathcal{S}^n \times \mathcal{Y}^{i-1} \rightarrow \mathcal{X} \quad (162)$$

$$(m, \mathbf{s}, y^{i-1}) \mapsto f_i(m, \theta^*, \mathbf{s}, y^{i-1}), \quad i \in [1 : n]. \quad (163)$$

It then follows from (160) that for every $m \in \mathcal{M}$ and $\mathbf{s} \in \mathcal{S}^n$

$$\sum_{\mathbf{y} \in \mathcal{Y}^n \setminus \mathcal{D}_m} \prod_{i=1}^n W(y_i | g_i(m, \mathbf{s}, y^{i-1}), s_i) = 0, \quad (164)$$

so the encoding mappings $\{g_i\}_{i \in [1:n]}$ and the decoding sets (159) constitute an (n, \mathcal{M}) zero-error feedback code with acausal SI and a deterministic encoder (6). \square

To conclude this section, we show that allowing for any (not necessarily discrete) random variable Θ in Definition A.1 does not lead to a more general notion of an (n, \mathcal{M}) zero-error feedback code with acausal SI and a stochastic encoder. To this end we shall use the following lemma, which is proved, e.g., in [16]:

Lemma A.2 (Functional Representation Lemma). *Given two chance variables X and Y of finite support, there exist a chance variable S of finite support \mathcal{S} that is independent of X and a function $g: \mathcal{X} \times \mathcal{S} \rightarrow \mathcal{Y}$ such that $Y = g(X, S)$.*

Remark A.3. *An (n, \mathcal{M}) zero-error feedback code with acausal SI and a stochastic encoder can also be viewed as a collection of n conditional PMFs*

$$P_{X_i | M, S^n, X^{i-1}, Y^{i-1}}, \quad i \in [1 : n] \quad (165)$$

and $|\mathcal{M}|$ disjoint decoding sets (159) for which (160) holds, where

$$\begin{aligned} & \mathbb{P}[Y^n \notin \mathcal{D}_m | M = m, S^n = \mathbf{s}] \\ &= \sum_{\mathbf{y} \in \mathcal{Y}^n \setminus \mathcal{D}_m} \sum_{\mathbf{x} \in \mathcal{X}^n} \prod_{i=1}^n P(x_i | m, \mathbf{s}, x^{i-1}, y^{i-1}) W(y_i | x_i, s_i). \end{aligned} \quad (166)$$

Indeed, for every (not necessarily discrete) random variable Θ of support \mathcal{T} , encoding mappings (158), and decoding sets (159), there exist n conditional PMFs (165) for which

$$\begin{aligned} & \sum_{\mathbf{y} \in \mathcal{Y}^n \setminus \mathcal{D}_m} \sum_{\mathbf{x} \in \mathcal{X}^n} \prod_{i=1}^n P(x_i | m, \mathbf{s}, x^{i-1}, y^{i-1}) W(y_i | x_i, s_i) \\ &= \mathbb{E}_{\Theta} \left[\sum_{\mathbf{y} \in \mathcal{Y}^n \setminus \mathcal{D}_m} \prod_{i=1}^n W(y_i | f_i(m, \Theta, \mathbf{s}, y^{i-1}), s_i) \right], \quad \forall m \in \mathcal{M}, \mathbf{s} \in \mathcal{S}^n. \end{aligned} \quad (167)$$

Conversely, for every collection of conditional PMFs (165) and decoding sets (159), there exist a random variable Θ of support \mathcal{T} and encoding mappings (158) for which (167) holds. Since the channel-input, -state, and -output alphabets are finite, repeated application of the Functional Representation lemma, Lemma A.2, yields, moreover, that we can choose the support \mathcal{T} of Θ finite as in Definition A.1.

B A Proof of Remarks 2.5 and 2.6

Proof. We begin with Remark 2.6. We first show that Condition (14) implies that the RHS of (11) is positive. To this end assume that (14) holds, pick $\mathcal{U} = \mathcal{Y}$, and let U be independent of \mathcal{S} and uniform over \mathcal{U} , so

$$P_{U|S} = P_U = \text{Unif}(\mathcal{U}). \quad (168)$$

Fix some conditional PMF $P_{X|U,S}$ that satisfies

$$\left(\left(W(u|x, s) > 0 \right) \implies \left(P_{X|U,S}(x|u, s) = 0 \right) \right), \quad \forall (u, x, s) \in \mathcal{U} \times \mathcal{X} \times \mathcal{S}. \quad (169)$$

(Such a $P_{X|U,S}$ exists, because (14) says that for every pair $(u, s) \in \mathcal{U} \times \mathcal{S}$ there exists some $\tilde{x} = \tilde{x}(u, s) \in \mathcal{X}$ for which $W(u|\tilde{x}, s)$ is zero, and we can thus choose $P_{X|U,S}$ to assign $\tilde{x}(u, s)$ probability one.) For $P_{U,X|S} = P_U \times P_{X|U,S}$, for every PMF P_S on \mathcal{S} , and for every conditional PMF $P_{Y|U,X,S}$ satisfying

$$P_{Y|U=u,X,S} \in \mathcal{P}(W), \quad \forall u \in \mathcal{U}, \quad (170)$$

we obtain w.r.t. the joint PMF $P_S \times P_{U,X|S} \times P_{Y|U,X,S}$

$$I(U; Y) - I(U; S) \stackrel{(a)}{=} I(U; Y) \quad (171)$$

$$\stackrel{(b)}{=} \log |\mathcal{Y}| - H(U|Y) \quad (172)$$

$$\stackrel{(c)}{\geq} \log |\mathcal{Y}| - \log(|\mathcal{Y}| - 1) \quad (173)$$

$$> 0, \quad (174)$$

where (a) holds because U is independent of S ; (b) holds because U is uniform over its support \mathcal{Y} ; and (c) holds because (169) and (170) imply that $(P_S \times P_{X|U,S} \times P_{Y|U,X,S})$ -almost-surely $U \neq Y$, and because the uniform distribution maximizes entropy. From this we conclude that Condition (14) is sufficient for the RHS of (11) to be positive.

We next turn to proving that if the RHS of (11) is positive, then (14) holds. We prove the contrapositive: we show that if for some $(s^*, y^*) \in \mathcal{S} \times \mathcal{Y}$

$$W(y^*|x, s^*) > 0, \forall x \in \mathcal{X}, \quad (175)$$

then the RHS of (11) must be zero. Suppose s^* and y^* are as above, introduce the PMF on \mathcal{S}

$$P_S(s) = \begin{cases} 1 & \text{if } s = s^*, \\ 0 & \text{otherwise,} \end{cases} \quad (176)$$

and choose $P_{Y|U,X,S} = P_{Y|X,S}$, where

$$P_{Y|X,S}(y|x, s) = \begin{cases} 1 & \text{if } s = s^*, y = y^*, \\ 0 & \text{if } s = s^*, y \neq y^*, \\ W(y|x, s) & \text{otherwise.} \end{cases} \quad (177)$$

Note that the conditional PMF $P_{Y|U,X,S}$ satisfies $P_{Y|U=u,X,S} \in \mathcal{P}(W)$, $\forall u \in \mathcal{U}$, because (175) and (177) imply that $P_{Y|X,S} \in \mathcal{P}(W)$. For every conditional PMF $P_{U,X|S}$, (176) and (177) imply that $(P_S \times P_{U,X|S} \times P_{Y|U,X,S})$ -almost-surely $Y = y^*$, and hence we obtain w.r.t. the joint PMF $P_S \times P_{U,X|S} \times P_{Y|U,X,S}$

$$I(U; Y) - I(U; S) \leq 0. \quad (178)$$

Since this holds for every conditional PMF $P_{U,X|S}$, we conclude that

$$\min_{P_S} \max_{P_{U,X|S}} \min_{\substack{P_{Y|U,X,S} \\ P_{Y|U=u,X,S} \in \mathcal{P}(W), \forall u \in \mathcal{U}}} I(U; Y) - I(U; S) = 0, \quad (179)$$

where the mutual informations are computed w.r.t. the joint PMF $P_S \times P_{U,X|S} \times P_{Y|U,X,S}$.

Having established Remark 2.6, we next prove Remark 2.5 by providing an example for which Theorem 2.3 implies that $C_{f,0} = 0$, and yet (14) holds. Such an example is the SD-DMC $W(y|x, s)$ for which $\mathcal{X} = \mathcal{Y} = \{0, 1, 2\}$ and

$$W(y|x, s) = \begin{cases} \frac{1}{2} & \text{if } y \neq x \oplus_3 2, \\ 0 & \text{otherwise.} \end{cases} \quad (180)$$

□

C Analysis of Example 2.9 where $W(y|x, s)$ is $\{0, 1\}$ -valued

In this appendix we assume that $W(y|x, s)$ is $\{0, 1\}$ -valued, and we derive (15) from Theorems 2.3 and 2.4.

We first show that Theorem 2.3 implies that $C_{f,0}$ is positive iff the RHS of (15) is positive. The latter is positive iff

$$|\{y \in \mathcal{Y}: \exists x \in \mathcal{X} \text{ s.t. } W(y|x, s) > 0\}| \geq 2, \forall s \in \mathcal{S}, \quad (181)$$

i.e., iff for every state there exists a pair of inputs that the deterministic channel maps to different outputs. By Theorem 2.3 $C_{f,0}$ is positive iff (10) holds, and we thus have to show that

$$(10) \iff (181). \quad (182)$$

The assumption that $W(y|x, s)$ is $\{0, 1\}$ -valued implies that for every pair of states $s, s' \in \mathcal{S}$ (not necessarily distinct) and every pair of inputs $x, x' \in \mathcal{X}$

$$W(y|x, s) W(y|x', s') = \begin{cases} 1 & \text{if } W(y|x, s) = W(y|x', s') = 1, \\ 0 & \text{otherwise.} \end{cases} \quad (183)$$

Using this we prove (182), beginning with

$$(10) \implies (181). \quad (184)$$

If we let $s' = s$, then (10) and (183) imply that for every state $s \in \mathcal{S}$ there exists a pair of inputs $x, x' \in \mathcal{X}$ that the channel maps to different outputs $y, y' \in \mathcal{Y}$, so

$$y \neq y' \quad \text{and} \quad W(y|x, s) = W(y'|x', s) = 1, \quad (185)$$

and hence

$$|\{y \in \mathcal{Y}: \exists x \in \mathcal{X} \text{ s.t. } W(y|x, s) > 0\}| \geq 2. \quad (186)$$

This proves (184). It remains to show that

$$(10) \longleftarrow (181). \quad (187)$$

From (181) and (183) it follows that for every state $s \in \mathcal{S}$ there exists a pair of inputs $x, x' \in \mathcal{X}$ that the deterministic channel maps to different outputs $y, y' \in \mathcal{Y}$, so

$$y \neq y' \quad \text{and} \quad W(y|x, s) = W(y'|x', s) = 1. \quad (188)$$

This implies that for every pair of states $s, s' \in \mathcal{S}$ (not necessarily distinct) there exists a pair of inputs $x, x' \in \mathcal{X}$ that the deterministic channel maps to different outputs $y, y' \in \mathcal{Y}$, so

$$y \neq y' \quad \text{and} \quad W(y|x, s) = W(y'|x', s') = 1, \quad (189)$$

and hence we conclude that (187) holds.

It remains to show that when $C_{f,0}$ is positive, then the RHS of (11) coincides with the RHS of (15). We first show that

$$C_{f,0} = \min_{P_S} \max_{P_{U,X|S}} I(U; Y) - I(U; S), \quad (190)$$

where the mutual informations are computed w.r.t. the joint PMF $P_S \times P_{U,X|S} \times W$. Note that for every $u \in \mathcal{U}$ the condition that $P_{Y|U=u, X, S} \in \mathcal{P}(W)$ is satisfied iff for every pair $(x, s) \in \mathcal{X} \times \mathcal{S}$ the outputs that have probability zero w.r.t. $W(\cdot|x, s)$ have probability zero w.r.t. $P_{Y|U, X, S}(\cdot|u, x, s)$. By the assumption that $W(y|x, s)$ is $\{0, 1\}$ -valued, this holds iff

$$P_{Y|U=u, X, S} = W, \quad \forall u \in \mathcal{U}, \quad (191)$$

and therefore (190) follows from Theorem 2.4.

With (190) at hand, we are now ready to show that the RHS of (11) is upper-bounded by the RHS of (15): w.r.t. the joint PMF $P_S \times P_{U,X|S} \times W$

$$C_{f,0} \stackrel{(a)}{=} \min_{P_S} \max_{P_{U,X|S}} I(U; Y) - I(U; S) \quad (192)$$

$$\stackrel{(b)}{\leq} \min_{P_S} \max_{P_{U,X|S}} I(U; Y, S) - I(U; S) \quad (193)$$

$$\stackrel{(c)}{=} \min_{P_S} \max_{P_{U,X|S}} I(U; Y|S) \quad (194)$$

$$\stackrel{(d)}{\leq} \min_{P_S} \max_{P_{U,X|S}} H(Y|S) \quad (195)$$

$$\stackrel{(e)}{\leq} \min_{s \in \mathcal{S}} \log |\{y \in \mathcal{Y} : \exists x \in \mathcal{X} \text{ s.t. } W(y|x, s) > 0\}|, \quad (196)$$

where (a) holds by (190); (b) holds because conditioning cannot increase entropy; (c) follows from the chain rule; (d) holds because conditional entropy is nonnegative; and (e) holds because the uniform distribution maximizes entropy, and because we can choose P_S to assign probability one to some $s \in \mathcal{S}$ that minimizes

$$\log |\{y \in \mathcal{Y} : \exists x \in \mathcal{X} \text{ s.t. } W(y|x, s) > 0\}|.$$

Having shown that the RHS of (11) is upper-bounded by the RHS of (15), we now conclude by showing that the reverse also holds, i.e., that the RHS of (11) is lower-bounded by the RHS of (15). Take $\mathcal{U} = \mathcal{Y}$, and for every $s \in \mathcal{S}$ choose $P_{U|S}(\cdot|s)$ to be the uniform distribution on the set

$$\{y \in \mathcal{Y} : \exists x \in \mathcal{X} \text{ s.t. } W(y|x, s) > 0\}.$$

By the assumption that $W(y|x, s)$ is $\{0, 1\}$ -valued, this choice of $P_{U|S}$ guarantees that for every pair $(u, s) \in \mathcal{U} \times \mathcal{S}$ for which $P_{U|S}(u|s) > 0$ there exists some $x = x(u, s) \in \mathcal{X}$ for which $W(u|x, s) = 1$. Now choose $P_{X|U, S}$ to assign $x(u, s)$ probability one. For $P_{U, X|S} = P_{U|S} \times P_{X|U, S}$ and for every PMF P_S on \mathcal{S} , we obtain $(P_S \times P_{U|S} \times P_{X|U, S} \times W)$ -almost-surely $U = Y$ and w.r.t. the joint PMF $P_S \times P_{U|S} \times P_{X|U, S} \times W$

$$I(U; Y) - I(U; S) \stackrel{(a)}{=} H(U|S) \tag{197}$$

$$\stackrel{(b)}{=} \sum_{s \in \mathcal{S}} P_S(s) \log |\{y \in \mathcal{Y} : \exists x \in \mathcal{X} \text{ s.t. } W(y|x, s) > 0\}| \tag{198}$$

$$\stackrel{(c)}{\geq} \min_{s \in \mathcal{S}} \log |\{y \in \mathcal{Y} : \exists x \in \mathcal{X} \text{ s.t. } W(y|x, s) > 0\}|, \tag{199}$$

where (a) holds because $(P_S \times P_{U|S} \times P_{X|U, S} \times W)$ -almost-surely $U = Y$; (b) holds because $P_{U|S}(\cdot|s)$ is for every $s \in \mathcal{S}$ the uniform distribution on

$$\{y \in \mathcal{Y} : \exists x \in \mathcal{X} \text{ s.t. } W(y|x, s) > 0\};$$

and (c) holds because the minimum of

$$\log |\{y \in \mathcal{Y} : \exists x \in \mathcal{X} \text{ s.t. } W(y|x, s) > 0\}|$$

over $s \in \mathcal{S}$ cannot be larger than its weighted average over $s \in \mathcal{S}$ with weights $P_S(s)$, $s \in \mathcal{S}$. From (190) and (199) we conclude that the RHS of (11) is lower-bounded by the RHS of (15).

D A Cardinality Bound on \mathcal{U}

Lemma D.1. *Given a channel $W(y|x, s)$ and a PMF P_S on \mathcal{S} , consider*

$$\max_{P_{U, X|S}} \min_{\substack{P_{Y|U, X, S}: \\ P_{Y|U=u, X, S} \in \mathcal{P}(W), \forall u \in \mathcal{U}}} I(U; Y) - I(U; S), \tag{200}$$

where the maximization is over all chance variables U of finite support, and the mutual informations are computed w.r.t. the joint PMF $P_S \times P_{U, X|S} \times P_{Y|U, X, S}$. Restricting X to be a function of U and S , i.e., $P_{U, X|S}$ to have the form

$$P_{U, X|S}(u, x|s) = P_{U|S}(u|s) \mathbb{1}_{x=g(u, s)}, \tag{201}$$

does not change (200). Nor does requiring that U take values in a set \mathcal{U} whose cardinality $|\mathcal{U}|$ satisfies

$$|\mathcal{U}| \leq |\mathcal{X}|^{|\mathcal{S}|}. \quad (202)$$

Proof. We first show that restricting X to be a function of U and S does not change (200). By the Functional Representation lemma (Lemma A.2), for every conditional PMF $P_{U,X|S}$, there exists a chance variable V of finite support \mathcal{V} and a function $h: \mathcal{U} \times \mathcal{V} \times \mathcal{S} \rightarrow \mathcal{X}$ such that

$$P_{U,X|S}(u, x|s) = \sum_{v \in \mathcal{V}} P_{U|S}(u|s) P_V(v) \mathbb{1}_{x=h(u,v,s)}. \quad (203)$$

Consequently, (200) is equal to

$$\max_{P_V, h(\cdot), P_{U|S}} \min_{P_{Y|U,X,S}: P_{Y|U=u,X,S} \in \mathcal{P}(W), \forall u \in \mathcal{U}} I(U; Y) - I(U; S), \quad (204)$$

where the maximization is over all chance variables V of finite support \mathcal{V} , functions $h: \mathcal{U} \times \mathcal{V} \times \mathcal{S} \rightarrow \mathcal{X}$, and conditional PMFs over a finite set \mathcal{U} ; and where the mutual informations are computed w.r.t. the joint PMF $P_S \times P_V \times P_{U,X|V,S} \times P_{Y|U,X,S}$, where $P_{U,X|V,S}$ is

$$P_{U,X|V,S}(u, x|v, s) = P_{U|S}(u|s) \mathbb{1}_{x=h(u,v,s)}. \quad (205)$$

Fix some PMF P_V on \mathcal{V} , some function $h: \mathcal{U} \times \mathcal{V} \times \mathcal{S} \rightarrow \mathcal{X}$, and some conditional PMF $P_{U|S}$, and let $(U, V, X, S) \sim P_S \times P_V \times P_{U,X|V,S}$, where $P_{U,X|V,S}$ is given in (205). Let $\tilde{P}_{Y|U,V,X,S}$ be some conditional PMF satisfying

$$\tilde{P}_{Y|U=u,V=v,X,S} \in \mathcal{P}(W), \quad \forall (u, v) \in \mathcal{U} \times \mathcal{V}, \quad (206)$$

and note that this implies that

$$\tilde{P}_{Y|U=u,X,S} \in \mathcal{P}(W), \quad \forall u \in \mathcal{U}, \quad (207)$$

where

$$\tilde{P}_{Y|U,X,S}(y|u, x, s) = \sum_{v \in \mathcal{V}} \frac{P_{U,V,X,S}(u, v, x, s)}{\sum_{v' \in \mathcal{V}} P_{U,V,X,S}(u, v', x, s)} \tilde{P}_{Y|U,V,X,S}(y|u, v, x, s). \quad (208)$$

W.r.t. the joint PMF $P_{U,V,X,S} \times \tilde{P}_{Y|U,V,X,S}$,

$$I(U, V; Y) - I(U, V; S) \stackrel{(a)}{=} I(U; Y) - I(U; S) + I(V; Y|U) \quad (209)$$

$$\stackrel{(b)}{\geq} I(U; Y) - I(U; S), \quad (210)$$

where (a) follows from the chain rule and the independence of V and (U, S) under $P_{U,V,X,S}$ (205); and (b) holds because mutual information is nonnegative. Since $P_{U,X,S} \times \tilde{P}_{Y|U,X,S}$ is obtained from $P_{U,V,X,S} \times \tilde{P}_{Y|U,V,X,S}$ by integrating V out (208),

$$\begin{aligned} I(U; Y) - I(U; S) & \quad \text{w.r.t. } P_{U,V,X,S} \times \tilde{P}_{Y|U,V,X,S} \\ &= I(U; Y) - I(U; S) \quad \text{w.r.t. } P_{U,X,S} \times \tilde{P}_{Y|U,X,S} \end{aligned} \quad (211)$$

$$= I(U; Y) - I(U; S) \quad \text{w.r.t. } P_{U,V,X,S} \times \tilde{P}_{Y|U,X,S}. \quad (212)$$

This and (210) imply that

$$\begin{aligned} I(U, V; Y) - I(U, V; S) & \quad \text{w.r.t. } P_{U,V,X,S} \times \tilde{P}_{Y|U,V,X,S} \\ & \geq I(U; Y) - I(U; S) \quad \text{w.r.t. } P_{U,V,X,S} \times \tilde{P}_{Y|U,X,S}. \end{aligned} \quad (213)$$

Since (206) implies (207), it follows from (213) that

$$\begin{aligned} & \max_{P_V, h(\cdot), P_{U|S}} \min_{\substack{\tilde{P}_{Y|U,V,X,S} \\ \tilde{P}_{Y|U=u, V=v, X, S} \in \mathcal{P}(W), \forall (u,v) \in \mathcal{U} \times \mathcal{V}}} I(U, V; Y) - I(U, V; S) \\ & \geq \max_{P_V, h(\cdot), P_{U|S}} \min_{\substack{\tilde{P}_{Y|U,X,S} \\ \tilde{P}_{Y|U=u, X, S} \in \mathcal{P}(W), \forall u \in \mathcal{U}}} I(U; Y) - I(U; S), \end{aligned} \quad (214)$$

where the mutual informations are computed w.r.t. $P_S \times P_V \times P_{U,X|V,S} \times \tilde{P}_{Y|U,V,X,S}$ in the first line and w.r.t. $P_S \times P_V \times P_{U,X|V,S} \times \tilde{P}_{Y|U,X,S}$ in the second line, where $P_{U,X|V,S}$ is given in (205). The RHS of (214) is (204), which, as we have noted, is equal to (200). Consequently, the LHS of (214) upper-bounds (200). But the LHS of (214) corresponds to choosing the auxiliary chance variable $\tilde{U} = (U, V)$, with the result that X is a deterministic function of (\tilde{U}, S) .

It remains to show that restricting the cardinality of \mathcal{U} to (202) does not change (200) when the maximization in (200) is over all conditional PMFs $P_{U,X|S}$ of the form (201). To this end we show that (200) does not change when we require that for every distinct $u_1, u_2 \in \mathcal{U}$ the mappings $g(u_1, \cdot)$ and $g(u_2, \cdot)$ differ. Since there are $|\mathcal{X}|^{|\mathcal{S}|}$ different mappings with domain \mathcal{S} and co-domain \mathcal{X} , this implies that restricting the cardinality of \mathcal{U} to (202) does not change (200).

Fix some finite set \mathcal{U} and some conditional PMF $P_{U,X|S}$ of the form (201), and let $(U, X, S) \sim P_S \times P_{U,X|S}$. Suppose that there exist distinct $u_1, u_2 \in \mathcal{U}$ for which

$$g(u_1, s) = g(u_2, s), \quad \forall s \in \mathcal{S}. \quad (215)$$

Define the chance variable

$$T = \begin{cases} U & \text{if } U \neq u_2, \\ u_1 & \text{otherwise} \end{cases} \quad (216)$$

of support $\mathcal{T} = \mathcal{U} \setminus \{u_2\}$, and denote by $P_{U,T,X,S}$ the joint PMF of (U, T, X, S) . By (215)

$$P_{X|T,S}(x|t, s) = \mathbb{1}_{x=g(t,s)}, \quad (217)$$

where

$$P_{X|T,S}(x|t, s) = \frac{\sum_{u \in \mathcal{U}} P_{U,T,X,S}(u, t, x, s)}{\sum_{(u', x') \in \mathcal{U} \times \mathcal{X}} P_{U,T,X,S}(u', t, x', s)}. \quad (218)$$

We will show that replacing U with T does not decrease our payoff, i.e., that

$$\begin{aligned} & \min_{\substack{\tilde{P}_{Y|U,X,S}: \\ \tilde{P}_{Y|U=u,X,S} \in \mathcal{P}(W), \forall u \in \mathcal{U}}} I(U; Y) - I(U; S) \\ & \leq \min_{\substack{\tilde{P}_{Y|T,X,S}: \\ \tilde{P}_{Y|T=t,X,S} \in \mathcal{P}(W), \forall t \in \mathcal{T}}} I(T; Y) - I(T; S), \end{aligned} \quad (219)$$

where the mutual informations are computed w.r.t. $P_{U,X,S} \times \tilde{P}_{Y|U,X,S}$ in the first line and w.r.t. $P_{T,X,S} \times \tilde{P}_{Y|T,X,S}$ in the second line. By repeating this process we can repeatedly reduce the cardinality of the support set of the auxiliary chance variable until $u_1 \neq u_2$ implies that $g(u_1, \cdot)$ and $g(u_2, \cdot)$ differ.

Let $\tilde{P}_{Y|T,X,S}$ be some conditional PMF satisfying

$$\tilde{P}_{Y|T=t,X,S} \in \mathcal{P}(W), \quad \forall t \in \mathcal{T}, \quad (220)$$

and define the conditional PMF

$$\tilde{P}_{Y|U,X,S}(y|u, x, s) = \begin{cases} \tilde{P}_{Y|T,X,S}(y|u, x, s) & \text{if } u \neq u_2, \\ \tilde{P}_{Y|T,X,S}(y|u_1, x, s) & \text{otherwise,} \end{cases} \quad (221)$$

so $\tilde{P}_{Y|U,X,S}(y|u, x, s) = \tilde{P}_{Y|T,X,S}(y|t, x, s)$ when $u = t$ or when $u = u_2$ and $t = u_1$. From this and (216), which implies that $P_{U,T,X,S}(u, t, x, s)$ is positive only when $u = t$ or when $u = u_2$ and $t = u_1$, it follows that

$$P_{U,T,X,S} \times \tilde{P}_{Y|U,X,S} = P_{U,T,X,S} \times \tilde{P}_{Y|T,X,S}. \quad (222)$$

From (220) and the definition of $\tilde{P}_{Y|U,X,S}$ (221) we see that

$$\tilde{P}_{Y|U=u,X,S} \in \mathcal{P}(W), \quad \forall u \in \mathcal{U}. \quad (223)$$

W.r.t. the joint PMF $P_{U,T,X,S} \times \tilde{P}_{Y|U,X,S}$ (which equals $P_{U,T,X,S} \times \tilde{P}_{Y|T,X,S}$ by (222))

$$\begin{aligned} I(U; Y) - I(U; S) & \stackrel{(a)}{=} I(T, U; Y) - I(T, U; S) \end{aligned} \quad (224)$$

$$\stackrel{(b)}{=} I(T; Y) - I(T; S) + I(U; Y|T) - I(U; S|T) \quad (225)$$

$$\stackrel{(c)}{=} I(T; Y) - I(T; S) + H(U|T, S) - H(U|T, Y), \quad (226)$$

where (a) holds because under $P_{U,T,X,S}$ T is determined by U (216); (b) follows from the chain rule; and (c) holds by definition of mutual information. W.r.t. the joint PMF $P_{U,T,X,S} \times \tilde{P}_{Y|U,X,S} = P_{U,T,X,S} \times \tilde{P}_{Y|T,X,S}$ the term $H(U|T, S) - H(U|T, Y)$ is not positive, because

$$\begin{aligned} H(U|T, S) - H(U|T, Y) & \stackrel{(a)}{=} H(U|T, X, S) - H(U|T, Y) \end{aligned} \quad (227)$$

$$\stackrel{(b)}{\leq} I(Y; U|T, X, S) \quad (228)$$

$$\stackrel{(c)}{=} 0, \quad (229)$$

where (a) holds because under $P_{U,T,X,S}$ X is determined by (T, S) (217); (b) holds because conditioning cannot increase entropy and by definition of mutual information; and (c) holds because under $P_{U,T,X,S} \times \tilde{P}_{Y|U,X,S} = P_{U,T,X,S} \times \tilde{P}_{Y|T,X,S}$ U and Y are conditionally independent given (T, X, S) . From (226) and (229) we obtain

$$\begin{aligned} I(U; Y) - I(U; S) & \quad \text{w.r.t. } P_{U,T,X,S} \times \tilde{P}_{Y|U,X,S} \\ & \leq I(T; Y) - I(T; S) \quad \text{w.r.t. } P_{U,T,X,S} \times \tilde{P}_{Y|T,X,S}, \end{aligned} \quad (230)$$

which is equivalent to

$$\begin{aligned} I(U; Y) - I(U; S) & \quad \text{w.r.t. } P_{U,X,S} \times \tilde{P}_{Y|U,X,S} \\ & \leq I(T; Y) - I(T; S) \quad \text{w.r.t. } P_{T,X,S} \times \tilde{P}_{Y|T,X,S}. \end{aligned} \quad (231)$$

Since (220) and (221) imply (223), we obtain from (231) that (219) holds, i.e., that replacing U with T does not decrease our payoff.

We can repeat the above process until we are left with a chance variable \bar{U} of finite support $\bar{\mathcal{U}} \subseteq \mathcal{U}$ that satisfies that for every distinct $\bar{u}_1, \bar{u}_2 \in \bar{\mathcal{U}}$ the mappings $g(\bar{u}_1, \cdot)$ and $g(\bar{u}_2, \cdot)$ differ, and, by (219), that

$$\begin{aligned} & \min_{\tilde{P}_{Y|U,X,S}:} I(U; Y) - I(U; S) \\ & \quad \tilde{P}_{Y|U=u,X,S} \in \mathcal{P}(W), \forall u \in \mathcal{U} \\ & \leq \min_{\tilde{P}_{Y|\bar{U},X,S}:} I(\bar{U}; Y) - I(\bar{U}; S). \\ & \quad \tilde{P}_{Y|\bar{U}=\bar{u},X,S} \in \mathcal{P}(W), \forall \bar{u} \in \bar{\mathcal{U}} \end{aligned} \quad (232)$$

From (232) we obtain the claim that (200)—with the maximization being over all conditional PMFs $P_{U,X|S}$ of the form (201)—does not change when we require that for every distinct $u_1, u_2 \in \mathcal{U}$ the mappings $g(u_1, \cdot)$ and $g(u_2, \cdot)$ differ. \square

E A Proof of Lemma 3.2

Proof. Recall that $\Theta = \lceil 2^{k(H(U|S)-\epsilon)} \rceil$. If $\epsilon \geq H(U|S)$, then $\Theta = 1$. The only size-1 partition of $\mathcal{T}_{P_U}^{(k)}$ is $\mathcal{T}_{P_U}^{(k)}$ itself, and, because $\mathcal{T}_{P_{U,S}}^{(k)}$ is not empty (since $P_{U,S}$ is a k -type), this partition satisfies (102), i.e.,

$$\forall \mathbf{s} \in \mathcal{T}_{P_S}^{(k)} \quad \exists \mathbf{u} \in \mathcal{T}_{P_U}^{(k)} \text{ s.t. } (\mathbf{u}, \mathbf{s}) \in \mathcal{T}_{P_{U,S}}^{(k)}. \quad (233)$$

Consider now the more interesting case where $\epsilon < H(U|S)$. We will show that if k exceeds some $\eta_0(|\mathcal{U}|, |\mathcal{S}|, \epsilon)$ (to be specified later), then the desired partition $\{\mathcal{B}_\ell\}_{\ell \in [1:\Theta]}$ of $\mathcal{T}_{P_U}^{(k)}$ exists. We shall do so using the probabilistic method. Fix $k \in \mathbb{N}$ and a k -type $P_{U,S}$ with corresponding conditional entropy $H(U|S)$. Generate a random partition $\{\mathcal{B}_\ell\}_{\ell \in [1:\Theta]}$ of $\mathcal{T}_{P_U}^{(k)}$, where $\{\mathcal{B}_\ell\}$ is short for $\{\mathcal{B}_\ell\}_{\ell \in [1:\Theta]}$, by placing each k -tuple $\mathbf{u} \in \mathcal{T}_{P_U}^{(k)}$ in a uniformly-drawn bin. We show that the probability that $\{\mathcal{B}_\ell\}$ violates (102) is smaller than one whenever $k \geq \eta_0(|\mathcal{U}|, |\mathcal{S}|, \epsilon)$. From this it will follow that the desired partition exists.

To upper-bound the probability that $\{\mathcal{B}_\ell\}$ violates (102), we first upper-bound

$$\mathbb{P} \left[\nexists \mathbf{u} \in \mathcal{B}_\ell \text{ s.t. } (\mathbf{u}, \mathbf{s}) \in \mathcal{T}_{P_{U,S}}^{(k)} \right]$$

for any fixed pair $(\mathbf{s}, \ell) \in \mathcal{T}_{P_S}^{(k)} \times [1 : \Theta]$:

$$\begin{aligned} & \mathbb{P} \left[\nexists \mathbf{u} \in \mathcal{B}_\ell \text{ s.t. } (\mathbf{u}, \mathbf{s}) \in \mathcal{T}_{P_{U,S}}^{(k)} \right] \\ &= \mathbb{P} \left[\mathcal{B}_\ell \cap \mathcal{T}_{P_{U|S}}^{(k)}(\mathbf{s}) = \emptyset \right] \end{aligned} \quad (234)$$

$$\stackrel{(a)}{=} (1 - \Theta^{-1})^{|\mathcal{T}_{P_{U|S}}^{(k)}(\mathbf{s})|} \quad (235)$$

$$\stackrel{(b)}{\leq} \left(1 - 2^{-k(H(U|S)-\epsilon)+1} \right)^{|\mathcal{T}_{P_{U|S}}^{(k)}(\mathbf{s})|} \quad (236)$$

$$\stackrel{(c)}{\leq} \exp \left\{ -2^{k\epsilon - \log(1+k)|\mathcal{U}||\mathcal{S}|+1} \right\}, \quad (237)$$

where (a) holds because each k -tuple $\mathbf{u} \in \mathcal{T}_{P_{U|S}}^{(k)}(\mathbf{s})$ is placed in \mathcal{B}_ℓ with probability Θ^{-1} ; (b) holds because

$$\Theta = \lceil 2^{k(H(U|S)-\epsilon)} \rceil \leq 2^{k(H(U|S)-\epsilon)+1}, \quad (238)$$

where the last inequality holds by assumption that $\epsilon < H(U|S)$; and (c) holds because $1 - \xi \leq e^{-\xi}$, $\xi \in \mathbb{R}$, and because $|\mathcal{T}_{P_{U|S}}^{(k)}(\mathbf{s})| \geq (1+k)^{-|\mathcal{U}||\mathcal{S}|} 2^{kH(U|S)}$ [8, Lemma 2.5]. Having obtained (237) for every fixed $(\mathbf{s}, \ell) \in \mathcal{T}_{P_S}^{(k)} \times [1 : \Theta]$, we use the Union-of-Events bound to upper-bound the probability that $\{\mathcal{B}_\ell\}$ violates (102):

$$\begin{aligned} & \mathbb{P}\left[\exists (\mathbf{s}, \ell) \in \mathcal{T}_{P_S}^{(k)} \times [1 : \Theta] \text{ s.t. } \mathcal{B}_\ell \cap \mathcal{T}_{P_{U|S}}^{(k)} = \emptyset\right] \\ & \stackrel{(a)}{\leq} |\mathcal{T}_{P_S}^{(k)}| \Theta \exp\left\{-2^{k\epsilon - \log(1+k)|\mathcal{U}||\mathcal{S}|+1}\right\} \end{aligned} \quad (239)$$

$$\stackrel{(b)}{\leq} \exp\left\{-2^{k\epsilon - \log(1+k)|\mathcal{U}||\mathcal{S}|+1} + k(\ln|\mathcal{S}| + \ln|\mathcal{U}| - \epsilon \ln 2)\right\}, \quad (240)$$

where (a) follows from the Union-of-Events bound and (237); and (b) holds because $|\mathcal{T}_{P_S}^{(k)}| \leq |\mathcal{S}|^k$ and $\Theta \leq |\mathcal{U}|^k$. The exponent on the RHS of (240),

$$-2^{k\epsilon - \log(1+k)|\mathcal{U}||\mathcal{S}|+1} + k(\ln|\mathcal{S}| + \ln|\mathcal{U}| - \epsilon \ln 2),$$

depends only on k , $|\mathcal{U}|$, $|\mathcal{S}|$, and ϵ , and it tends to $-\infty$ as k tends to infinity. Consequently, there exists some $\eta_0(|\mathcal{U}|, |\mathcal{S}|, \epsilon)$ that guarantees that the exponent is negative whenever $k \geq \eta_0(|\mathcal{U}|, |\mathcal{S}|, \epsilon)$. For such values of k the RHS of (240) is smaller than one, and the desired partition exists. \square

F A Proof of Theorem 2.10

The proof consists of a direct and a converse part. We first establish the direct part.

Direct Part. If there exists a partition $\mathcal{Y}_0, \mathcal{Y}_1$ of \mathcal{Y} satisfying (16), then the encoder can transmit a bit $m \in \{0, 1\}$ error-free in one channel use: If $m = 0$ and the Time-1 channel-state is $s \in \mathcal{S}$, then it sends some $x \in \mathcal{X}$ for which $W(\mathcal{Y}_0|x, s) = 1$, and if $m = 1$ and the Time-1 channel-state is $s \in \mathcal{S}$, then it sends some $x' \in \mathcal{X}$ for which $W(\mathcal{Y}_1|x', s) = 1$. This allows the decoder to recover the transmitted bit error-free by declaring “ $m = 0$ ” if the Time-1 channel-output is in \mathcal{Y}_0 and “ $m = 1$ ” if the Time-1 channel-output is in \mathcal{Y}_1 . \square

We next prove the converse part of Theorem 2.10.

Converse Part. To prove that (16) is necessary for $C_{f,0}^{\text{caus}}$ to be positive, we will show that if no partition $\mathcal{Y}_0, \mathcal{Y}_1$ of \mathcal{Y} satisfies (16), then it is impossible to transmit a bit error-free. Assume then that no such partition exists, and let the bit take values in the set $\mathcal{M} = \{0, 1\}$. Fix a blocklength n and n encoding mappings

$$f_i: \mathcal{M} \times \mathcal{S}^i \times \mathcal{Y}^{i-1} \rightarrow \mathcal{X}, \quad i \in [1 : n].$$

To show that the mappings do not achieve error-free transmission, we will exhibit a pair of state sequences $\mathbf{s}, \tilde{\mathbf{s}} \in \mathcal{S}^n$ and an output sequence $\mathbf{y} \in \mathcal{Y}^n$ that for every $i \in [1 : n]$ satisfy

$$W(y_i | f_i(0, s^i, y^{i-1}), s_i) W(y_i | f_i(1, \tilde{s}^i, y^{i-1}), \tilde{s}_i) > 0. \quad (241)$$

This will rule out error-free transmission, because if the state sequence is either \mathbf{s} or $\tilde{\mathbf{s}}$, then the decoder, not knowing which, cannot recover the bit.

Our construction of $\mathbf{s}, \tilde{\mathbf{s}} \in \mathcal{S}^n$ and $\mathbf{y} \in \mathcal{Y}^n$ is inductive, i.e., we first exhibit Time-1 components $s_1, \tilde{s}_1 \in \mathcal{S}$ and $y_1 \in \mathcal{Y}$ that satisfy (241) for $i = 1$, and we then repeatedly increment i by one (until it reaches n) and exhibit Time- i components $s_i, \tilde{s}_i \in \mathcal{S}$ and $y_i \in \mathcal{Y}$ that—together with the previously constructed $\{s_j, \tilde{s}_j\}_{j \in [1:i-1]}$ and $\{y_j\}_{j \in [1:i-1]}$ —satisfy (241).

We start by exhibiting Time-1 components $s_1, \tilde{s}_1 \in \mathcal{S}$ and $y_1 \in \mathcal{Y}$ that satisfy (241) for $i = 1$. To this end we show that

$$\exists s, \tilde{s} \in \mathcal{S}, y \in \mathcal{Y} \text{ s.t. } W(y | f_1(0, s), s) W(y | f_1(1, \tilde{s}), \tilde{s}) > 0. \quad (242)$$

Our proof of (242) is by contradiction. To reach a contradiction, suppose that (242) does not hold, so

$$\left(W(y | f_1(0, s), s) W(y | f_1(1, \tilde{s}), \tilde{s}) = 0, \forall y \in \mathcal{Y} \right), \forall s, \tilde{s} \in \mathcal{S}. \quad (243)$$

Define the set

$$\mathcal{Y}_0 = \left\{ y \in \mathcal{Y} : \exists s \in \mathcal{S} \text{ s.t. } W(y | f_1(0, s), s) > 0 \right\} \quad (244)$$

and its \mathcal{Y} -complement $\mathcal{Y}_1 = \mathcal{Y} \setminus \mathcal{Y}_0$. By the definition of the set \mathcal{Y}_0

$$W(\mathcal{Y}_0 | f_1(0, s), s) = 1, \forall s \in \mathcal{S}, \quad (245)$$

and by (243)

$$W(\mathcal{Y}_0 | f_1(1, \tilde{s}), \tilde{s}) = 0, \forall \tilde{s} \in \mathcal{S}, \quad (246)$$

so

$$W(\mathcal{Y}_1 | f_1(1, s), s) = 1 - W(\mathcal{Y}_0 | f_1(1, s), s) = 1, \forall s \in \mathcal{S}. \quad (247)$$

This contradicts our assumption that no partition $\mathcal{Y}_0, \mathcal{Y}_1$ of \mathcal{Y} satisfies (16) and thus establishes (242). If s, \tilde{s} , and y are as promised in (242), then we choose $s_1 = s, \tilde{s}_1 = \tilde{s}$, and $y_1 = y$ with the result that (241) holds for $i = 1$.

For the inductive step, suppose $\ell \in [2 : n]$, and that we have already constructed $\{s_i, \tilde{s}_i\}_{i \in [1:\ell-1]}$ and $\{y_i\}_{i \in [1:\ell-1]}$ for which (241) holds for every $i \in [1 : \ell - 1]$. We construct Time- ℓ components $s_\ell, \tilde{s}_\ell \in \mathcal{S}$ and $y_\ell \in \mathcal{Y}$ that—together with the previously constructed $\{s_i, \tilde{s}_i\}_{i \in [1:\ell-1]}$ and $\{y_i\}_{i \in [1:\ell-1]}$ —satisfy (241) when we substitute ℓ for i in (241), i.e., we show that

$$\exists s_\ell, \tilde{s}_\ell \in \mathcal{S}, y_\ell \in \mathcal{Y} \text{ s.t. } W(y_\ell | f_\ell(0, s^\ell, y^{\ell-1}), s_\ell) W(y_\ell | f_\ell(1, \tilde{s}^\ell, y^{\ell-1}), \tilde{s}_\ell) > 0. \quad (248)$$

Our proof of (248) is by contradiction. To reach a contradiction, suppose that (248) does not hold, so

$$\left(W(y_\ell | f_\ell(0, s^\ell, y^{\ell-1}), s_\ell) W(y_\ell | f_\ell(1, \tilde{s}^\ell, y^{\ell-1}), \tilde{s}_\ell) = 0, \forall y_\ell \in \mathcal{Y} \right), \forall s_\ell, \tilde{s}_\ell \in \mathcal{S}. \quad (249)$$

Define the set

$$\mathcal{Y}_0 = \left\{ y_\ell \in \mathcal{Y} : \exists s_\ell \in \mathcal{S} \text{ s.t. } W(y_\ell | f_\ell(0, s^\ell, y^{\ell-1}), s_\ell) > 0 \right\} \quad (250)$$

and its \mathcal{Y} -complement $\mathcal{Y}_1 = \mathcal{Y} \setminus \mathcal{Y}_0$. By the definition of the set \mathcal{Y}_0

$$W(\mathcal{Y}_0 | f_\ell(0, s^\ell, y^{\ell-1}), s_\ell) = 1, \forall s_\ell \in \mathcal{S}, \quad (251)$$

and by (249)

$$W(\mathcal{Y}_0 | f_\ell(1, \tilde{s}^\ell, y^{\ell-1}), \tilde{s}_\ell) = 0, \forall \tilde{s}_\ell \in \mathcal{S}, \quad (252)$$

so

$$W(\mathcal{Y}_1 | f_\ell(1, \tilde{s}^\ell, y^{\ell-1}), \tilde{s}_\ell) = 1 - W(\mathcal{Y}_0 | f_\ell(1, \tilde{s}^\ell, y^{\ell-1}), \tilde{s}_\ell) = 1, \forall \tilde{s}_\ell \in \mathcal{S}. \quad (253)$$

This contradicts our assumption that no partition $\mathcal{Y}_0, \mathcal{Y}_1$ of \mathcal{Y} satisfies (16).

Since the construction goes through for every $\ell \in [1 : n]$, when ℓ reaches n we have constructed a pair of state sequences $\mathbf{s}, \tilde{\mathbf{s}} \in \mathcal{S}^n$ and an output sequence $\mathbf{y} \in \mathcal{Y}^n$ that for every $i \in [1 : n]$ satisfy (241). \square

G A Proof of Theorem 2.11

Suppose $W(y|x, s)$ satisfies the condition in Theorem 2.10 for $C_{f,0}^{\text{caus}}$ to be positive. In this case the RHS of (17) and the RHS of (18) are equal, because the latter is the zero-error feedback capacity of the (state-less) DMC $W'(y|u)$ (3) and thus—using Ahlswede's alternative form (4)—can be alternatively expressed as (17). It thus suffices to prove (18), i.e.,

$$C_{f,0}^{\text{caus}} = \max_{P_U} \min_y -\log \sum_{u: W'(y|u) > 0} P_U(u). \quad (254)$$

The proof consists of a direct and a converse part. We first establish the direct part.

Direct Part. That the RHS of (254) is achievable follows from Shannon's results on the zero-error capacity [1, Theorem 7] and on channels with states [11]. Indeed, the encoder can convert the channel to a state-less channel whose inputs are Shannon strategies [11]. That is, it can perform the encoding over the set \mathcal{U} , where $\{g(u, \cdot) : u \in \mathcal{U}\}$ equals $\mathcal{X}^{\mathcal{S}}$, and transmit at Time i the channel input $g(u_i(m), S_i)$, where $u_i(m)$ is the i -th component of the codeword $\mathbf{u}(m)$ corresponding to the message m to be transmitted (see Figure 3 and [12, Remark 7.6]). In doing so, the encoder transforms the SD-DMC $W(y|x, s)$ with causal SI and feedback into the state-less DMC

$$W'(y|u) = \sum_{s \in \mathcal{S}} Q_S(s) W(y|g(u, s), s)$$

with feedback. Because the zero-error feedback capacity of the DMC $W'(y|u)$ is equal to the RHS of (254) (see [1, Theorem 7] or (3)), the RHS of (254) is achievable. \square

We next establish the converse part.

Converse Part. To establish that $C_{f,0}^{\text{caus}}$ cannot be larger than the RHS of (254), we adapt Shannon's converse of [1, Theorem 7] to the present setting. Let

$$\xi = \max_{P_U} \min_y -\log \sum_{u: W'(y|u) > 0} P_U(u), \quad (255)$$

and fix a finite set \mathcal{M} , a blocklength n , and n encoding mappings

$$f_i: \mathcal{M} \times \mathcal{S}^i \times \mathcal{Y}^{i-1} \rightarrow \mathcal{X}, \quad i \in [1 : n].$$

We will exhibit an output sequence $\mathbf{y} \in \mathcal{Y}^n$ for which the corresponding post- n survivor-set

$$\mathcal{M}_n = \left\{ m \in \mathcal{M} : \exists \mathbf{s} \in \mathcal{S}^n \text{ s.t. } \prod_{i=1}^n W(y_i | f_i(m, \mathbf{s}^i, \mathbf{y}^{i-1}), s_i) > 0 \right\} \quad (256)$$

is of size at least

$$|\mathcal{M}_n| \geq 2^{-n\xi} |\mathcal{M}|. \quad (257)$$

From (256) and (257) it will then follow that the probability of a decoding error can only be zero if $|\mathcal{M}| \leq 2^{n\xi}$, because otherwise $|\mathcal{M}_n| \geq 2$ and none of the messages in \mathcal{M}_n can be ruled out by the decoder.

To conclude the proof, we show by mathematical induction over $i \in [0 : n]$ that for every $i \in [0 : n]$

$$\exists \mathbf{y} \in \mathcal{Y}^i \text{ s.t. } (|\mathcal{M}_i(\mathbf{y})| \geq 2^{-i\xi} |\mathcal{M}|), \quad (258a)$$

where $\mathcal{M}_i(\mathbf{y})$ is the post- i survivor-set corresponding to \mathbf{y} , so

$$\mathcal{M}_i(\mathbf{y}) = \left\{ m \in \mathcal{M} : \exists \mathbf{s} \in \mathcal{S}^i \text{ s.t. } \prod_{j=1}^i W(y_j | f_j(m, s^j, y^{j-1}), s_j) > 0 \right\}. \quad (258b)$$

In (258b) we use the convention that the empty product is 1, so $\mathcal{M}_0(\emptyset) = \mathcal{M}$ for $i = 0$.

1. Basis $i = 0$: Because $\mathcal{M}_0(\emptyset) = \mathcal{M}$, (258) holds for $i = 0$.
2. Inductive Step: Fix $\ell \in [1 : n]$, and assume that (258) holds for $i = \ell - 1$, i.e., that there exists some $y^{\ell-1} \in \mathcal{Y}^{\ell-1}$ for which

$$|\mathcal{M}_{\ell-1}(y^{\ell-1})| \geq 2^{-(\ell-1)\xi} |\mathcal{M}|. \quad (259)$$

Suppose $y^{\ell-1}$ is as above. By the definition of the set $\mathcal{M}_{\ell-1}(y^{\ell-1})$ (258b) there exists a collection $\{s^{\ell-1}(m)\}_{m \in \mathcal{M}_{\ell-1}(y^{\ell-1})}$ of $(\ell - 1)$ -tuples from $\mathcal{S}^{\ell-1}$ for which

$$\prod_{i=1}^{\ell-1} W(y_i | f_i(m, s^i(m), y^{i-1}), s_i(m)) > 0, \quad \forall m \in \mathcal{M}_{\ell-1}(y^{\ell-1}). \quad (260)$$

To prove that (258a) holds for $i = \ell$, we show that

$$\begin{aligned} \exists y \in \mathcal{Y} \text{ s.t. } \left(\left| \left\{ m \in \mathcal{M}_{\ell-1}(y^{\ell-1}) : \exists s_\ell(m) \in \mathcal{S} \text{ s.t. } \right. \right. \right. \\ \left. \left. \left. W(y | f_\ell(m, s^\ell(m), y^{\ell-1}), s_\ell(m)) > 0 \right\} \right| \right. \\ \left. \geq 2^{-\xi} |\mathcal{M}_{\ell-1}(y^{\ell-1})| \right). \end{aligned} \quad (261)$$

Setting y_ℓ to be the $y \in \mathcal{Y}$ promised in (261) will prove (258) for $i = \ell$.

Because $\{g(u, \cdot) : u \in \mathcal{U}\}$ equals $\mathcal{X}^\mathcal{S}$, for every $m \in \mathcal{M}_{\ell-1}(y^{\ell-1})$ there exists a $u \in \mathcal{U}$, call it $u_\ell(m)$, satisfying that

$$f_\ell(m, s^\ell(m), y^{\ell-1}) = g(u, s_\ell(m)), \quad \forall s_\ell(m) \in \mathcal{S}. \quad (262)$$

This and (20) imply that (261) is equivalent to

$$\begin{aligned} \exists y \in \mathcal{Y} \text{ s.t. } \left(\left| \left\{ m \in \mathcal{M}_{\ell-1}(y^{\ell-1}) : W'(y | u_\ell(m)) > 0 \right\} \right| \right. \\ \left. \geq 2^{-\xi} |\mathcal{M}_{\ell-1}(y^{\ell-1})| \right). \end{aligned} \quad (263)$$

It thus suffices to establish (263). The proof is essentially the converse of [1, Theorem 7]. For every $u \in \mathcal{U}$ denote by F_u the fraction of all the messages $m \in \mathcal{M}_{\ell-1}(y^{\ell-1})$ for which $u_\ell(m)$ equals u , so

$$F_u \triangleq \frac{|\{m \in \mathcal{M}_{\ell-1}(y^{\ell-1}) : u_\ell(m) = u\}|}{|\mathcal{M}_{\ell-1}(y^{\ell-1})|}, \quad u \in \mathcal{U}. \quad (264)$$

The construction of the collection $\{F_u\}_{u \in \mathcal{U}}$ guarantees that for every $y \in \mathcal{Y}$

$$\left| \left\{ m \in \mathcal{M}_{\ell-1}(y^{\ell-1}) : W'(y|u_\ell(m)) > 0 \right\} \right| = \sum_{u: W'(y|u) > 0} F_u |\mathcal{M}_{\ell-1}(y^{\ell-1})|. \quad (265)$$

Moreover, the collection $\{F_u\}_{u \in \mathcal{U}}$ is like a PMF on \mathcal{U} , i.e.,

$$F_u \geq 0, \quad \forall u \in \mathcal{U}, \quad (266a)$$

$$\sum_{u \in \mathcal{U}} F_u = 1. \quad (266b)$$

Choose y as one that—among all elements of \mathcal{Y} —maximizes

$$\sum_{u: W'(y|u) > 0} F_u. \quad (267)$$

For this choice of y we obtain the lower bound

$$\begin{aligned} & \left| \left\{ m \in \mathcal{M}_{\ell-1}(y^{\ell-1}) : W'(y|u_\ell(m)) > 0 \right\} \right| \\ & \stackrel{(a)}{=} \sum_{u: W'(y|u) > 0} F_u |\mathcal{M}_{\ell-1}(y^{\ell-1})| \end{aligned} \quad (268)$$

$$\stackrel{(b)}{=} \max_{y \in \mathcal{Y}} \sum_{u: W'(y|u) > 0} F_u |\mathcal{M}_{\ell-1}(y^{\ell-1})| \quad (269)$$

$$\stackrel{(c)}{\geq} \min_{P_U} \max_{y \in \mathcal{Y}} \sum_{u: W'(y|u) > 0} P_U(u) |\mathcal{M}_{\ell-1}(y^{\ell-1})| \quad (270)$$

$$\stackrel{(d)}{\geq} 2^{-\xi} |\mathcal{M}_{\ell-1}(y^{\ell-1})|, \quad (271)$$

where (a) holds by (265); (b) holds because y maximizes (267) and consequently also (265) among all elements of \mathcal{Y} ; (c) holds by (266); and (d) holds by (255). This proves (263) and consequently also (261). If y is as promised in (261) and we choose y_ℓ to be y , then it follows from (259) and (260) that for $i = \ell$ the post- ℓ survivor-set $\mathcal{M}_\ell(y^\ell)$ of (258b) is of size at least $2^{-\ell\xi} |\mathcal{M}|$, and hence that (258) holds for $i = \ell$. □

H A Proof of Remarks 2.12 and 2.13

Proof. We begin with Remark 2.13. We first show that Condition (22) implies that the RHS of (17) is positive. To this end assume that (22) holds. Recall that $\{g(u, \cdot) : u \in \mathcal{U}\}$

equals $\mathcal{X}^{\mathcal{S}}$. This, combined with (22), implies that for every $y \in \mathcal{Y}$ there must exist a $u \in \mathcal{U}$, call it u_y , that satisfies

$$W(y|g(u_y, s), s) = 0, \quad \forall s \in \mathcal{S}. \quad (272)$$

The mapping $y \mapsto u_y$ need not be one-to-one, but it follows from (272) that the cardinality of its range must exceed one. Let U be uniform over the set $\{u_y : y \in \mathcal{Y}\}$, so

$$P_U = \text{Unif}(\{u_y : y \in \mathcal{Y}\}). \quad (273)$$

For every $P_{Y|U} \in \mathcal{P}(W')$ we obtain w.r.t. $P_U \times P_{Y|U}$

$$I(U; Y) \stackrel{(a)}{=} \log |\{u_y : y \in \mathcal{Y}\}| - H(U|Y) \quad (274)$$

$$\stackrel{(b)}{\geq} \log |\{u_y : y \in \mathcal{Y}\}| - \log(|\{u_y : y \in \mathcal{Y}\}| - 1) \quad (275)$$

$$\stackrel{(c)}{\geq} \log |\mathcal{Y}| - \log(|\mathcal{Y}| - 1), \quad (276)$$

where (a) holds because U is uniform over $\{u_y : y \in \mathcal{Y}\}$; (b) holds because $U \neq u_Y$ and because the uniform distribution maximizes entropy; and (c) holds because $|\mathcal{Y}| \geq 2$ (which follows from (22)), and because the function

$$\xi \mapsto \frac{\xi}{\xi - 1}, \quad \xi > 1$$

is strictly monotonically decreasing in ξ . This implies that the RHS of (17) is positive:

$$\begin{aligned} \max_{P_U} \min_{P_{Y|U} \in \mathcal{P}(W')} I(U; Y) \\ \geq \log |\mathcal{Y}| - \log(|\mathcal{Y}| - 1) \end{aligned} \quad (277)$$

$$> 0, \quad (278)$$

where the mutual information is computed w.r.t. the joint PMF $P_U \times P_{Y|U}$.

We next turn to proving that if the RHS of (17) is positive, then (22) holds. We prove the contrapositive: we show that if for some $(s^*, y^*) \in \mathcal{S} \times \mathcal{Y}$

$$W(y^*|x, s^*) > 0, \quad \forall x \in \mathcal{X}, \quad (279)$$

then the RHS of (17) must be zero. Suppose s^* and y^* are as above, and introduce the conditional PMF

$$P_{Y|U}(y|u) = \begin{cases} 1 & \text{if } y = y^*, \\ 0 & \text{otherwise.} \end{cases} \quad (280)$$

Note that $P_{Y|U} \in \mathcal{P}(W')$, because (279) implies that

$$W(y^*|g(u, s^*), s^*) > 0, \forall u \in \mathcal{U}. \quad (281)$$

For every PMF P_U on \mathcal{U} (280) implies that $(P_U \times P_{Y|U})$ -almost-surely $Y = y^*$, and hence we obtain w.r.t. the joint PMF $P_U \times P_{Y|U}$

$$I(U; Y) = 0. \quad (282)$$

Because this holds for every PMF P_U on \mathcal{U} , we conclude that the RHS of (17) is zero:

$$\max_{P_U} \min_{P_{Y|U} \in \mathcal{P}(W')} I(U; Y) = 0, \quad (283)$$

where the mutual information is computed w.r.t. the joint PMF $P_U \times P_{Y|U}$.

Having established Remark 2.13, we next prove Remark 2.12 by providing an example for which Theorem 2.10 implies that $C_{f,0}^{\text{caus}} = 0$, and yet (22) holds. Such an example is the SD-DMC $W(y|x, s)$ for which $\mathcal{X} = \mathcal{Y} = \{0, 1, 2\}$ and

$$W(y|x, s) = \begin{cases} \frac{1}{2} & \text{if } y \neq x \oplus_3 2, \\ 0 & \text{otherwise.} \end{cases} \quad (284)$$

□

I A Proof of Theorem 2.19

The proof consists of a direct and a converse part. We first establish the direct part.

Direct Part. We assume that (10) holds and show that the RHS of (26) is achievable. If the RHS of (26) is zero, then there is nothing to prove, so we assume that it is positive. The proof builds on the proofs of Remark 3.1 and the direct part of Theorem 2.4, adapting both to the case where—in addition to the message—the encoder wants to convey to the receiver error-free also the state sequence. We partition the blocklength- n transmission into $B + 1$ blocks, with each of the first B blocks being of length k , and with Block $(B + 1)$ being of length k' . The choice we shall later make for k' will be such that the last block be of negligible length compared to Bk and therefore not affect the code's asymptotic rate.

Before the transmission begins, the encoder is revealed the realization $\mathbf{s} \triangleq S^n$ of the state sequence, from which it can compute the realization $\mathbf{s}^{(b)} \triangleq s_{(b-1)k+1}^{bk}$ of the

Block- b state-sequence for every $b \in [1 : B]$ and the realization $\mathbf{s}^{(B+1)} \triangleq s_{Bk+1}^{Bk+k'}$ of the Block- $(B+1)$ state-sequence. In the first B blocks our scheme draws on the scheme we used in the direct part of Theorem 2.4 but with the following two modifications: 1) to guarantee that the decoder can recover the Block- $(B+1)$ state-sequence $\mathbf{s}^{(B+1)}$, the encoder transmits the pair $(m, \mathbf{s}^{(B+1)}) \in \mathcal{M} \times \mathcal{S}^{k'}$ comprising the message to be sent and the Block- $(B+1)$ state-sequence; and 2) to guarantee that the decoder can recover the state sequences $\{\mathbf{s}^{(b)}\}_{b \in [1:B]}$ during the first B blocks, we choose the auxiliary chance variable U to comprise the channel state S and consequently to be (X, S) (because we can w.l.g. restrict X to be a function of U and S). The last block draws on Phase 2 of the scheme we used to prove Remark 3.1. We next describe the proposed coding scheme in detail, beginning with the first B blocks and ending with the last block.

For every $b \in [1 : B]$ we adapt the Block b transmission of the scheme we used in the direct part of Theorem 2.4 as follows. Assume for now that the decoder—while incognizant of $\mathbf{s}^{(1)}, \dots, \mathbf{s}^{(B)}$ —knows the empirical types $P_{\mathbf{s}^{(1)}}, \dots, P_{\mathbf{s}^{(B)}}$: Block $(B+1)$ will ensure that the scheme works even though the decoder is incognizant of these types. Let $\mathcal{I}_0 = \mathcal{M} \times \mathcal{S}^{k'}$ be the set of all possible pairs of message $m' \in \mathcal{M}$ and Block- $(B+1)$ state-sequence $\mathbf{s}' \in \mathcal{S}^{k'}$, and for every $b \in [1 : B]$ let $\mathcal{I}_b \subset \mathcal{M} \times \mathcal{S}^{bk} \times \mathcal{S}^{k'}$ be the (random) set comprising all the triples of message $m' \in \mathcal{M}$, state sequence $\hat{\mathbf{s}} \in \mathcal{S}^{bk}$ pertaining to the first b blocks, and Block- $(B+1)$ state-sequence $\mathbf{s}' \in \mathcal{S}^{k'}$ that have a positive posterior probability given the channel outputs Y^{bk} and the empirical types $\{P_{\mathbf{s}^{(b')}}\}_{b' \in [1:b]}$ during the first b blocks. Choose some k -type $P_{X,S}^{(b)}$ whose \mathcal{S} -marginal $P_S^{(b)}$ equals $P_{\mathbf{s}^{(b)}}$. In the following, unless otherwise specified, all entropies and mutual informations are computed w.r.t. the joint PMF $P_{X,S}^{(b)}$.

For any k -length state-sequence $\mathbf{s}' \in \mathcal{T}_{P_S^{(b)}}^{(k)}$, let $L_{P_{X|S}^{(b)}}^{(k)}$ denote the size of the $P_{X|S}^{(b)}(\mathbf{s}')$ -shell $\mathcal{T}_{P_{X|S}^{(b)}}^{(k)}(\mathbf{s}')$, i.e.,

$$L_{P_{X|S}^{(b)}}^{(k)} = \left| \mathcal{T}_{P_{X|S}^{(b)}}^{(k)}(\mathbf{s}') \right|, \quad \mathbf{s}' \in \mathcal{T}_{P_S^{(b)}}^{(k)}. \quad (285)$$

This size does not depend on $\mathbf{s}' \in \mathcal{T}_{P_S^{(b)}}^{(k)}$, and by [8, Lemma 2.5]

$$L_{P_{X|S}^{(b)}}^{(k)} \geq (1+k)^{-|\mathcal{X}||\mathcal{S}|} 2^{kH(X|S)}. \quad (286)$$

We partition \mathcal{I}_{b-1} into $L_{P_{X|S}^{(b)}}^{(k)}$ subsets whose size is between

$$\left\lfloor |\mathcal{I}_{b-1}| / L_{P_{X|S}^{(b)}}^{(k)} \right\rfloor \quad \text{and} \quad \left\lceil |\mathcal{I}_{b-1}| / L_{P_{X|S}^{(b)}}^{(k)} \right\rceil;$$

and we associate with each set a different bin from the bins

$$\mathcal{B}_\ell \subseteq \mathcal{T}_{P_{X,S}^{(b)}}^{(k)}, \quad \ell \in \left[1 : L_{P_{X|S}^{(b)}}^{(k)}\right],$$

where the bins $\{\mathcal{B}_\ell\}$ are pairwise disjoint subsets of $\mathcal{T}_{P_{X,S}^{(b)}}^{(k)}$

$$\mathcal{B}_\ell \cap \mathcal{B}_{\ell'} = \emptyset, \quad \left(\forall \ell, \ell' \in \left[1 : L_{P_{X|S}^{(b)}}^{(k)}\right], \ell' \neq \ell \right), \quad (287a)$$

and where each bin “covers” $\mathcal{T}_{P_S^{(b)}}^{(k)}$ exactly in the sense that

$$\forall (\mathbf{s}, \ell) \in \mathcal{T}_{P_S^{(b)}}^{(k)} \times \left[1 : L_{P_{X|S}^{(b)}}^{(k)}\right] \quad \exists! \mathbf{x} \in \mathcal{T}_{P_X^{(b)}}^{(k)} \text{ s.t. } (\mathbf{x}, \mathbf{s}) \in \mathcal{B}_\ell. \quad (287b)$$

(Unlike the direct part of Theorem 2.4, here we need not invoke Lemma 3.2 to guarantee the existence of such bins. Indeed, that such bins exist follows from the definition of $L_{P_{X|S}^{(b)}}^{(k)}$ (285): for every $\mathbf{s} \in \mathcal{T}_{P_S^{(b)}}^{(k)}$ there exist $L_{P_{X|S}^{(b)}}^{(k)}$ different \mathbf{x} for which $(\mathbf{x}, \mathbf{s}) \in \mathcal{T}_{P_{X,S}^{(b)}}^{(k)}$ (285), and hence we can choose some collection $\{\mathcal{B}_\ell\}$ satisfying that for every $\mathbf{s} \in \mathcal{T}_{P_S^{(b)}}^{(k)}$ each of the $L_{P_{X|S}^{(b)}}^{(k)}$ pairs in $\mathcal{T}_{P_{X,S}^{(b)}}^{(k)}$ whose second component is \mathbf{s} is contained in a different bin from the $L_{P_{X|S}^{(b)}}^{(k)}$ bins $\{\mathcal{B}_\ell\}$.) To transmit the triple $(m, s^{bk}, \mathbf{s}^{(B+1)})$, the encoder picks from the bin that is associated with the subset of \mathcal{I}_{b-1} containing $(m, s^{(b-1)k}, \mathbf{s}^{(B+1)})$ the pair $(\mathbf{x}', \mathbf{s}')$ satisfying $\mathbf{s}' = \mathbf{s}^{(b)}$ (287b) and chooses as the Block- b channel-inputs $\mathbf{x}^{(b)}$ the k -tuple \mathbf{x}' .

Based on the Block- b outputs $\mathbf{y}^{(b)} \triangleq Y_{(b-1)k+1}^{bk}$ and the empirical type $P_{\mathbf{s}^{(b)}}$, the encoder and decoder compute \mathcal{I}_b as follows. First, they identify all the pairs $(\tilde{\mathbf{x}}, \tilde{\mathbf{s}}) \in \mathcal{T}_{P_{X,S}^{(b)}}^{(k)}$ that could have produced the observed Block- b outputs $\mathbf{y}^{(b)}$. For each such pair $(\tilde{\mathbf{x}}, \tilde{\mathbf{s}})$ they identify the unique bin that contains it, and they include in \mathcal{I}_b all the triples $(m', \hat{\mathbf{s}}, \mathbf{s}') \in \mathcal{M} \times \mathcal{S}^{bk} \times \mathcal{S}^{k'}$ satisfying that $\hat{s}_{(b-1)k+1}^{bk} = \tilde{\mathbf{s}}$ and that $(m', \hat{s}^{(b-1)k}, \mathbf{s}')$ is an element of the subset of \mathcal{I}_{b-1} with which this bin is associated.

Using arguments similar to those in the direct part of Theorem 2.4, we next show that

$$|\mathcal{I}_b| \leq \left(\max_{P_{Y|X,S} \in \mathcal{P}(W)} 2^{-k(I(X,S;Y) - H(S) - \beta_k)} \right) |\mathcal{I}_{b-1}|, \quad (288a)$$

whenever

$$|\mathcal{I}_{b-1}| \geq L_{P_{X|S}^{(b)}}^{(k)}, \quad (288b)$$

and

$$|\mathcal{I}_b| \leq 2^{k(H(X,S)+\beta_k)} \quad (288c)$$

otherwise, where the mutual information is computed w.r.t. the joint PMF $P_{X,S}^{(b)} \times P_{Y|X,S}$, and where β_k is given by

$$\beta_k = \frac{\log(1+k) |\mathcal{X}| |\mathcal{S}| (1 + |\mathcal{Y}|) + 1}{k} \quad (289)$$

and hence converges to zero as k tends to infinity. To this end note that, with probability one, the empirical type of the tuple $(\mathbf{x}^{(b)}, \mathbf{s}^{(b)}, \mathbf{y}^{(b)})$ satisfies

$$P_{\mathbf{x}^{(b)}, \mathbf{s}^{(b)}} = P_{X,S}^{(b)}, \quad (290a)$$

$$(W(y|x, s) = 0) \implies (P_{\mathbf{x}^{(b)}, \mathbf{s}^{(b)}, \mathbf{y}^{(b)}}(x, s, y) = 0). \quad (290b)$$

This allows us to upper-bound the number of pairs in $\mathcal{T}_{P_{X,S}^{(b)}}^{(k)}$ that could have produced the observed Block- b outputs $\mathbf{y}^{(b)}$: For every fixed k -type $P_{X,S,Y}$ on $\mathcal{X} \times \mathcal{S} \times \mathcal{Y}$ the number of pairs $(\tilde{\mathbf{x}}, \tilde{\mathbf{s}}) \in \mathcal{T}_{P_{X,S}^{(b)}}^{(k)}$ that satisfy $(\tilde{\mathbf{x}}, \tilde{\mathbf{s}}, \mathbf{y}^{(b)}) \in \mathcal{T}_{P_{X,S,Y}}^{(k)}$ cannot exceed $2^{kH(X,S|Y)}$, where the conditional entropy is computed w.r.t. the joint PMF $P_{X,S,Y}$ [8, Lemma 2.5]. This, combined with (290) and the fact that the number of k -types on $\mathcal{X} \times \mathcal{S} \times \mathcal{Y}$ cannot exceed $(1+k)^{|\mathcal{X}|+|\mathcal{S}|+|\mathcal{Y}|}$, implies that the number of pairs in $\mathcal{T}_{P_{X,S}^{(b)}}^{(k)}$ that could have produced the observed Block- b outputs $\mathbf{y}^{(b)}$ is upper-bounded by

$$2^{\log(1+k) |\mathcal{X}|+|\mathcal{S}|+|\mathcal{Y}|} \max_{P_{Y|X,S} \in \mathcal{P}(W)} 2^{kH(X,S|Y)}, \quad (291)$$

where the conditional entropy is computed w.r.t. the joint PMF $P_{X,S}^{(b)} \times P_{Y|X,S}$. Since the bins are pairwise disjoint (287a), no pair is contained in more than one bin. Every bin is associated with a subset of \mathcal{I}_{b-1} whose size is at most $\left\lceil |\mathcal{I}_{b-1}| / L_{P_{X|S}^{(b)}}^{(k)} \right\rceil$; and by (286)

$$\left\lceil |\mathcal{I}_{b-1}| / L_{P_{X|S}^{(b)}}^{(k)} \right\rceil \leq 2^{-kH(X|S)+\log(1+k) |\mathcal{X}|+|\mathcal{S}|+1} |\mathcal{I}_{b-1}|, \quad (292)$$

whenever (288b) holds, and

$$\left\lceil |\mathcal{I}_{b-1}| / L_{P_{X|S}^{(b)}}^{(k)} \right\rceil = 1 \quad (293)$$

otherwise. From (291)–(293), the fact that

$$H(X|S) - H(X, S|Y) = I(X, S; Y) - H(S), \quad (294)$$

and the inequality

$$H(X, S|Y) \leq H(X, S), \quad (295)$$

which holds because conditioning cannot increase entropy, we obtain (288).

We next use (288) to show that—for some choice of the k -type $P_{X,S}^{(b)}$ and some $\gamma_k = \gamma_k(|\mathcal{X}|, |\mathcal{S}|, |\mathcal{Y}|)$, which converges to zero as k tends to infinity—we can guarantee that

$$|\mathcal{I}_b| \leq \left(\max_{P_S} \min_{P_{X|S}} \max_{P_{Y|X,S} \in \mathcal{P}(W)} 2^{-k(I(X,S;Y) - H(S) - \gamma_k)} \right) |\mathcal{I}_{b-1}|, \quad (296a)$$

whenever

$$|\mathcal{I}_{b-1}| \geq 2^{k \log |\mathcal{X}|}, \quad (296b)$$

and

$$|\mathcal{I}_b| \leq 2^{k(\log |\mathcal{X}| + \log |\mathcal{S}| + \gamma_k)} \quad (296c)$$

otherwise, where the mutual information and the entropy are computed w.r.t. the joint PMF $P_S \times P_{X|S} \times P_{Y|X,S}$. To this end we will first infer from (288) that

$$|\mathcal{I}_b| \leq \left(\max_{P_{Y|X,S} \in \mathcal{P}(W)} 2^{-k(I(X,S;Y) - H(S) - \beta_k)} \right) |\mathcal{I}_{b-1}|, \quad (297a)$$

whenever

$$|\mathcal{I}_{b-1}| \geq 2^{k \log |\mathcal{X}|}, \quad (297b)$$

and

$$|\mathcal{I}_b| \leq 2^{k(\log |\mathcal{X}| + \log |\mathcal{S}| + \beta_k)} \quad (297c)$$

otherwise, where the mutual information is computed w.r.t. the joint PMF $P_{X,S}^{(b)} \times P_{Y|X,S}$, and where β_k is defined in (289). The following three observations show that (288) \implies (297) :

- 1) By (285) $L_{P_{X|S}^{(b)}}^{(k)} \leq |\mathcal{X}|^k$, so whenever Condition (297b) holds so does (288b). Consequently, we obtain from (288) that whenever Condition (297b) holds the inequality (297a) holds.
- 2) Since

$$\log |\mathcal{X}| - (I(X, S; Y) - H(S) - \beta_k) \leq \log |\mathcal{X}| + \log |\mathcal{S}| + \beta_k, \quad (298)$$

it follows from (288) that the inequality (297c) holds whenever

$$L_{P_{X|S}^{(b)}}^{(k)} \leq |\mathcal{I}_{b-1}| < 2^{k \log |\mathcal{X}|}. \quad (299)$$

- 3) Since $H(X, S) \leq \log |\mathcal{X}| + \log |\mathcal{S}|$, it follows from (288) that the inequality (297c) holds whenever

$$|\mathcal{I}_{b-1}| < L_{P_{X|S}^{(b)}}^{(k)}. \quad (300)$$

Having established (297), we are now ready to prove (296). Since we can choose any k -type $P_{X,S}^{(b)}$ whose \mathcal{S} -marginal $P_S^{(b)}$ is $P_{\mathbf{s}^{(b)}}$, we can choose $P_{X,S}^{(b)} = P_{\mathbf{s}^{(b)}} \times P_{X|S}^{(b)}$, where $P_{X|S}^{(b)}$ is the conditional k -type that—among all conditional k -types—maximizes

$$\min_{P_{Y|X,S} \in \mathcal{P}(W)} I(X, S; Y) - H(S), \quad (301)$$

where the mutual information and the entropy are computed w.r.t. the joint PMF $P_{X,S}^{(b)} \times P_{Y|X,S}$. Every conditional PMF can be approximated in the total variation distance by a conditional k -type when k is sufficiently large; and, because entropy and mutual information are continuous in this distance [8, Lemma 2.7], it follows that—for the above choice of the conditional k -type and some $\gamma_k = \gamma_k(|\mathcal{X}|, |\mathcal{S}|, |\mathcal{Y}|)$, which converges to zero as k tends to infinity—(297) implies (296).

Since we assume that the RHS of (26) is positive, we can choose B and k sufficiently large so that

$$\begin{aligned} & \left(\max_{P_S} \min_{P_{X|S}} \max_{P_{Y|X,S} \in \mathcal{P}(W)} 2^{-Bk(I(X,S;Y) - H(S) - \gamma_k)} \right) |\mathcal{M}| |\mathcal{S}|^{k'} \\ & \leq 2^{k(\log |\mathcal{X}| + \log |\mathcal{S}| + \gamma_k)}; \end{aligned} \quad (302)$$

and by (296) this guarantees that, with probability one,

$$|\mathcal{I}_B| \leq 2^{k(\log |\mathcal{X}| + \log |\mathcal{S}| + \gamma_k)}. \quad (303)$$

We now deal with Block $(B+1)$. Because the decoder is incognizant of the empirical types $\{P_{\mathbf{s}^{(b)}}\}_{b \in [1:B]}$, it cannot compute the post-Block- B ambiguity-set \mathcal{I}_B comprising the pairs of message and length- n state-sequence of positive posterior probability given the channel outputs $\{\mathbf{y}^{(b)}\}_{b \in [1:B]}$ and the k -types $\{P_{\mathbf{s}^{(b)}}\}_{b \in [1:B]}$. The uncertainty that needs to be addressed is about the message, the length- n state-sequence, as well as the B empirical types of $\mathbf{s}^{(1)}, \dots, \mathbf{s}^{(B)}$. Let $\mathcal{J}_B \subseteq \mathcal{M} \times \mathcal{S}^n$ denote the union of the post-Block- B ambiguity-sets corresponding to all the different B -tuples of k -types on \mathcal{S} , i.e., \mathcal{J}_B is the set of pairs of messages and state sequences that have a positive posterior probability given only the outputs $\{\mathbf{y}^{(b)}\}_{b \in [1:B]}$ (and not the k -types $\{P_{\mathbf{s}^{(b)}}\}_{b \in [1:B]}$). Because the post-Block- B ambiguity-set corresponding to any given B -tuple of k -types on \mathcal{S} satisfies (303), and because there are at most $(1+k)^B |\mathcal{S}|$ B -tuples of k -types on \mathcal{S} ,

$$|\mathcal{J}_B| \leq 2^{k(\log |\mathcal{X}| + \log |\mathcal{S}| + \gamma_k) + B \log(1+k) |\mathcal{S}|}. \quad (304)$$

In Block $(B+1)$ we resolve the set \mathcal{J}_B . This will guarantee that the decoder can recover the transmitted message m and the length- n state-sequence \mathbf{s} error-free.

Block $(B+1)$ is similar to Phase 2 of the scheme we used to prove Remark 3.1: the encoder allocates to every pair $(m', \mathbf{s}') \in \mathcal{J}_B$ a length- k' codeword $\mathbf{x}(m', \mathbf{s}')$, where the codewords are chosen so that

$$\begin{aligned} & \left(\forall (m', \mathbf{s}'), (m'', \mathbf{s}'') \in \mathcal{J}_B \text{ s.t. } m' \neq m'' \right) \quad \exists i \in [1 : k'] \text{ s.t.} \\ & \left(W(y | x_i(m', \mathbf{s}'), s'_{Bk+i}) W(y | x_i(m'', \mathbf{s}''), s''_{Bk+i}) = 0, \forall y \in \mathcal{Y} \right). \end{aligned} \quad (305)$$

(We shall shortly use a random coding argument to show that this can be done.) To convey the message m and the state sequence \mathbf{s} , the encoder transmits in Block $(B+1)$ the codeword $\mathbf{x}(m, \mathbf{s})$. Condition (305) implies that, upon observing the Block- $(B+1)$ outputs $\mathbf{y}^{(B+1)} \triangleq Y_{Bk+1}^{Bk+k'}$, the decoder, who knows \mathcal{J}_B and the codewords $\{\mathbf{x}(m', \mathbf{s}')\}$, can determine the transmitted message m and the state sequence \mathbf{s} error-free, because

$$\prod_{i=1}^{k'} W(y_i^{(B+1)} | x_i(m, \mathbf{s}), s_{Bk+i}) > 0, \quad (306)$$

whereas (305) implies for every other pair $(m', \mathbf{s}') \in \mathcal{J}_B$

$$\prod_{i=1}^{k'} W(y_i^{(B+1)} | x_i(m', \mathbf{s}'), s'_{Bk+i}) = 0. \quad (307)$$

The decoder can thus calculate $\prod_i W(y_i^{(B+1)} | x_i(\tilde{m}, \tilde{\mathbf{s}}), \tilde{s}_{Bk+i})$ for each $(\tilde{m}, \tilde{\mathbf{s}}) \in \mathcal{J}_B$ and produce the pair $(\tilde{m}, \tilde{\mathbf{s}})$ for which this product is positive.

We next show that, for some choice of k' , there exist codewords $\{\mathbf{x}(m', \mathbf{s}')\}$ satisfying (305). To this end we use a random coding argument. Draw the length- k' codewords $\{\mathbf{X}(m', \mathbf{s}')\}$ independently, each uniformly over $\mathcal{X}^{k'}$. From (10) it then follows that for any fixed distinct $(m', \mathbf{s}'), (m'', \mathbf{s}'') \in \mathcal{J}_B$

$$\begin{aligned} & \mathbb{P} \left[\forall i \in [1 : k'] \exists y \in \mathcal{Y} \text{ s.t. } W(y | X_i(m', \mathbf{s}'), s'_{Bk+i}) W(y | X_i(m'', \mathbf{s}''), s''_{Bk+i}) > 0 \right] \\ & \leq \left(1 - \frac{1}{|\mathcal{X}|^2} \right)^{k'} \end{aligned} \quad (308)$$

$$= 2^{-k'(2 \log |\mathcal{X}| - \log(|\mathcal{X}|^2 - 1))}. \quad (309)$$

This, the Union-of-Events bound, and (304) imply that the probability that the randomly drawn length- k' codewords do not satisfy (305) is upper-bounded by

$$\begin{aligned} & |\mathcal{J}_B|^2 2^{-k'(2 \log |\mathcal{X}| - \log(|\mathcal{X}|^2 - 1))} \\ & \leq 2^{-k'(2 \log |\mathcal{X}| - \log(|\mathcal{X}|^2 - 1)) + 2(k(\log |\mathcal{X}| + \log |\mathcal{S}| + \gamma_k) + B \log(1+k) |\mathcal{S}|)}, \end{aligned} \quad (310)$$

which is smaller than one whenever

$$k' > \frac{k(\log |\mathcal{X}| + \log |\mathcal{S}| + \gamma_k) + B \log(1+k) |\mathcal{S}|}{\log |\mathcal{X}| - \frac{1}{2} \log(|\mathcal{X}|^2 - 1)}. \quad (311)$$

Consequently, if we choose some k' that satisfies (311), then there exist length- k' code-words $\{\mathbf{x}(m', s')\}$ satisfying (305).

We are now ready to join the dots and conclude that the coding scheme asymptotically achieves any rate smaller than the RHS of (26). More precisely, we will show that, for every rate R smaller than the RHS of (26) and every sufficiently-large blocklength n , our coding scheme can convey nR bits and the length- n state-sequence error-free in n channel uses.

It follows from (302) and (311) that if the positive integers n , B , k , k' are such that (311) holds,

$$n = Bk + k', \quad (312)$$

and

$$nR + k' \log |\mathcal{S}| \leq Bk \left(\min_{P_S} \max_{P_{X|S}} \min_{P_{Y|X,S} \in \mathcal{P}(W)} I(X, S; Y) - H(S) - \gamma_k \right), \quad (313)$$

then our coding scheme can convey nR bits and the length- n state-sequence error-free in n channel uses. It thus remains to exhibit positive integers B , k , k' such that for every sufficiently-large blocklength n (311)–(313) hold. As we argue next, when n is sufficiently large we can choose

$$B = \lfloor \sqrt{n} \rfloor - \left(\left\lfloor \frac{\log |\mathcal{X}| + \log |\mathcal{S}| + \gamma_k + \log(1 + \sqrt{n}) |\mathcal{S}|}{\log |\mathcal{X}| - \frac{1}{2} \log(|\mathcal{X}|^2 - 1)} \right\rfloor + 1 \right), \quad (314a)$$

$$k = \lfloor \sqrt{n} \rfloor, \quad (314b)$$

$$k' = n - Bk. \quad (314c)$$

Note that, whenever n is sufficiently large, B , k , and k' are positive, and (311) and (312) are satisfied. To see that (313) holds whenever n is sufficiently large, we first observe from (314b) that k tends to infinity as n tends to infinity. Because $\gamma_k = \gamma_k(|\mathcal{X}|, |\mathcal{S}|, |\mathcal{Y}|)$ converges to zero as k tends to infinity, this implies that γ_k converges to zero as n tends to infinity. We next observe that (314) implies that Bk/n converges to one as n tends to infinity and consequently that $k'/n = 1 - Bk/n$ converges to zero as n tends to infinity. This, combined with the facts that γ_k converges to zero as n tends to infinity and that R is smaller than the RHS of (26), implies that (313) holds whenever n is sufficiently large. \square

We next establish the converse part of Theorem 2.19.

Converse Part. That (10) is a necessary condition for $C_{f,0}^{m+s}$ to be positive follows from Theorem 2.3, because $C_{f,0}^{m+s}$ is upper-bounded by $C_{f,0}$. We next show that—irrespective of whether or not (10) holds— $C_{f,0}^{m+s}$ is upper-bounded by the RHS of (26). The proof is similar to the converse of Theorem 2.4. Fix a finite set \mathcal{M} , a blocklength n , and an (n, \mathcal{M}) zero-error state-conveying code with n encoding mappings

$$f_i: \mathcal{M} \times \mathcal{S}^n \times \mathcal{Y}^{i-1} \rightarrow \mathcal{X}, \quad i \in [1 : n] \quad (315)$$

and $|\mathcal{M}| |\mathcal{S}|^n$ disjoint decoding sets $\mathcal{D}_{m,s} \subseteq \mathcal{Y}^n$, $(m, s) \in \mathcal{M} \times \mathcal{S}^n$. We will show that the rate $\frac{1}{n} \log |\mathcal{M}|$ of the code is upper-bounded by the RHS of (26).

Draw M uniformly over \mathcal{M} , and denote its distribution P_M . Since the code is a zero-error state-conveying code,

$$\mathbb{P}[Y^n \in \mathcal{D}_{M,S^n}] = 1, \quad (316)$$

where \mathbb{P} is the distribution (118) of (M, S^n, X^n, Y^n) induced by P_M , the state distribution Q , the encoding mappings (315), and the channel law $W(y|x, s)$. Similarly as in the converse of Theorem 2.4, fix any PMF \tilde{P}_S on \mathcal{S} and any collection of n conditional PMFs $\{\tilde{P}_{Y_i|X_i, S_i}\}_{i \in [1:n]}$ that satisfy

$$\tilde{P}_{Y_i|X_i, S_i} \in \mathcal{P}(W). \quad (317)$$

These PMFs induce the PMF on $\mathcal{M} \times \mathcal{S}^n \times \mathcal{X}^n \times \mathcal{Y}^n$

$$\tilde{P}_{M, S^n, X^n, Y^n} = P_M \times \tilde{P}_S^n \times \prod_{i=1}^n (P_{X_i|M, S^n, Y^{i-1}} \times \tilde{P}_{Y_i|X_i, S_i}). \quad (318)$$

It follows from (1) and (317) that $\tilde{P}_{M, S^n, X^n, Y^n} \ll \mathbb{P}$ and consequently that (316) implies

$$\tilde{P}_{M, S^n, X^n, Y^n}[Y^n \in \mathcal{D}_{M, S^n}] = 1. \quad (319)$$

We upper-bound $\frac{1}{n} \log |\mathcal{M}|$ by carrying out the following calculation under $\tilde{P}_{M, S^n, X^n, Y^n}$

of (319):

$$\begin{aligned} & \frac{1}{n} \log |\mathcal{M}| \\ & \stackrel{(a)}{=} \frac{1}{n} \left[H(M) + H(S^n) - H(S^n) \right] \end{aligned} \quad (320)$$

$$\stackrel{(b)}{=} \frac{1}{n} \left[I(S^n, M; Y^n) - H(S^n) \right] \quad (321)$$

$$\stackrel{(c)}{=} \frac{1}{n} \sum_{i=1}^n \left[I(S^n, M; Y_i | Y^{i-1}) - H(S_i | S^{i-1}) \right] \quad (322)$$

$$\stackrel{(d)}{\leq} \frac{1}{n} \sum_{i=1}^n \left[I(S^n, M, Y^{i-1}; Y_i) - H(S_i) \right] \quad (323)$$

$$\stackrel{(e)}{\leq} \frac{1}{n} \sum_{i=1}^n \left[I(X_i, S_i; Y_i) - H(S_i) \right], \quad (324)$$

where (a) holds because M is uniform over \mathcal{M} under $\tilde{P}_{M, S^n, X^n, Y^n}$; (b) holds by (319) and because M is independent of S^n under $\tilde{P}_{M, S^n, X^n, Y^n}$; (c) follows from the chain rule; (d) holds because conditioning cannot increase entropy and by the independence of S_i and S^{i-1} under $\tilde{P}_{M, S^n, X^n, Y^n}$; and (e) holds because under $\tilde{P}_{M, S^n, X^n, Y^n}$ (S^n, M, Y^{i-1}), (X_i, S_i) , and Y_i form a Markov chain in that order.

We will conclude the proof by exhibiting a PMF \tilde{P}_S and a collection of conditional PMFs $\{\tilde{P}_{Y_i | X_i, S_i}\}_{i \in [1:n]}$ satisfying (317) for which each summand on the RHS of (324) is upper-bounded by the RHS of (26).

We begin with the choice of $\{\tilde{P}_{Y_i | X_i, S_i}\}_{i \in [1:n]}$. We first choose $\tilde{P}_{Y_i | X_i, S_i}$ for $i = 1$, and we then repeatedly increment i by one until it reaches n . Key to our choice is the observation, which will be justified shortly, that \tilde{P}_{X_i, S_i} is determined by \tilde{P}_S and $\{\tilde{P}_{Y_j | X_j, S_j}\}_{j \in [i-1]}$. Our choice of $\tilde{P}_{Y_i | X_i, S_i}$ can thus depend not only on our choice of \tilde{P}_S and our previous choices of $\{\tilde{P}_{Y_j | X_j, S_j}\}_{j \in [1:i-1]}$ but also on \tilde{P}_{X_i, S_i} . This will allow us to choose $\tilde{P}_{Y_i | X_i, S_i}$ as one that—among all conditional PMFs satisfying (317)—minimizes

$$I(X_i, S_i; Y_i) - H(S_i), \quad (325)$$

where the mutual information and the entropy are computed w.r.t. the joint PMF $\tilde{P}_{X_i, S_i} \times \tilde{P}_{Y_i | X_i, S_i}$. Since (318) implies that

$$\tilde{P}_{S_i} = \tilde{P}_S, \quad i \in [1 : n], \quad (326)$$

we will then find that, for our choice of $\{\tilde{P}_{Y_i | X_i, S_i}\}$,

$$\begin{aligned} & I(X_i, S_i; Y_i) - H(S_i) \\ & \leq \max_{\tilde{P}_{X_i | S_i}} \min_{\tilde{P}_{Y_i | X_i, S_i} \in \mathcal{P}(W)} I(X_i, S_i; Y_i) - H(S_i), \quad i \in [1 : n], \end{aligned} \quad (327)$$

where the mutual information and the entropy are computed w.r.t. the joint PMF $\tilde{P}_{S_i} \times \tilde{P}_{X_i|S_i} \times \tilde{P}_{Y_i|X_i,S_i}$. The chosen conditional PMFs $\{\tilde{P}_{Y_i|X_i,S_i}\}_{i \in [1:n]}$ satisfy (317), and hence (324), (326), and (327) will imply that

$$\begin{aligned} & \frac{1}{n} \log |\mathcal{M}| \\ & \leq \max_{\tilde{P}_{X|S}} \min_{\tilde{P}_{Y|X,S} \in \mathcal{P}(W)} I(X, S; Y) - H(S), \end{aligned} \quad (328)$$

where the mutual information and the entropy in the i -th summand are computed w.r.t. the joint PMF $\tilde{P}_S \times \tilde{P}_{X|S} \times \tilde{P}_{Y|X,S}$.

We now prove that indeed \tilde{P}_{X_i,S_i} is determined by \tilde{P}_S and $\{\tilde{P}_{Y_j|X_j,S_j}\}_{j \in [1:i-1]}$. In fact, we will show that the latter two determine $\tilde{P}_{M,S^n,X^i,Y^{i-1}}$. The latter determines \tilde{P}_{X_i,S_i} , because the tuple (X_i, S_i) is determined by (M, S^n, X^i, Y^{i-1}) .

We use mathematical induction, but first we note that the PMF $\tilde{P}_{M,S^n,X^n,Y^n}$ is constructed inductively: by (318)

$$\tilde{P}_{M,S^n,X_1} = P_M \times \tilde{P}_S^n \times P_{X_1|M,S^n}, \quad (329)$$

and, for every $\ell \in [2 : n]$, $\tilde{P}_{M,S^n,X^\ell,Y^{\ell-1}}$ is constructed from $\tilde{P}_{M,S^n,X^{\ell-1},Y^{\ell-2}}$ by

$$\tilde{P}_{M,S^n,X^\ell,Y^{\ell-1}} = \tilde{P}_{M,S^n,X^{\ell-1},Y^{\ell-2}} \times \tilde{P}_{Y_{\ell-1}|X_{\ell-1},S_{\ell-1}} \times P_{X_\ell|M,S^n,Y^{\ell-1}}. \quad (330)$$

In describing the proof we shall make the dependence on P_M , our choice of \tilde{P}_S , and $\{P_{X_j|M,S^n,Y^{j-1}}\}_{j \in [1:n]}$, whose components are determined by the encoding mappings (315) via (119), implicit.

1. Basis $\ell = 1$: It follows from (329) that \tilde{P}_{M,S^n,X_1} is determined.
2. Inductive Step: Fix $\ell \in [2 : i]$, and suppose that $\tilde{P}_{M,S^n,X^{\ell-1},Y^{\ell-2}}$ is determined by $\{\tilde{P}_{Y_j|X_j,S_j}\}_{j \in [1:\ell-2]}$. This implies that $\tilde{P}_{M,S^n,X^{\ell-1},Y^{\ell-2}}$ and $\tilde{P}_{Y_{\ell-1}|X_{\ell-1},S_{\ell-1}}$ are determined by $\{\tilde{P}_{Y_j|X_j,S_j}\}_{j \in [1:\ell-1]}$. Consequently, it follows from (330) that $\tilde{P}_{M,S^n,X^\ell,Y^{\ell-1}}$ is determined by $\{\tilde{P}_{Y_j|X_j,S_j}\}_{j \in [1:\ell-1]}$.

This proves that, for every $i \in [1 : n]$, $\tilde{P}_{M,S^n,X^i,Y^{i-1}}$ and consequently also \tilde{P}_{X_i,S_i} are determined by \tilde{P}_S and $\{\tilde{P}_{Y_j|X_j,S_j}\}_{j \in [1:i-1]}$, and hence (328) holds.

Having established (328), we are now ready to conclude the proof. Since we can choose any PMF \tilde{P}_S on \mathcal{S} , we can choose one that—among all PMFs on \mathcal{S} —yields the tightest bound, i.e., minimizes

$$\max_{\tilde{P}_{X|S}} \min_{\tilde{P}_{Y|X,S} \in \mathcal{P}(W)} I(X, S; Y) - H(S), \quad (331)$$

where the mutual information and the entropy are computed w.r.t. the joint PMF $P_S \times P_{X|S} \times P_{Y|X,S}$. For this choice of \tilde{P}_S (328) implies that

$$\frac{1}{n} \log |\mathcal{M}| \leq \min_{\tilde{P}_S} \max_{P_{X|S}} \min_{\tilde{P}_{Y|X,S} \in \mathcal{P}(W)} I(X, S; Y) - H(S), \quad (332)$$

where the mutual information and the entropy are computed w.r.t. the joint PMF $P_S \times P_{X|S} \times P_{Y|X,S}$. If the RHS of (332) is negative, then—irrespective of $|\mathcal{M}| \geq 1$ —(332) is a contradiction and consequently (319) cannot hold. This implies that—even if $|\mathcal{M}|$ is one—the state sequence cannot be conveyed error-free. Since we say that $C_{f,0}^{m+s} = 0$ if the state sequence cannot be conveyed error-free, (332) implies that $C_{f,0}^{m+s}$ is upper-bounded by the RHS of (26). \square

J A Proof of Theorem 2.20

We already showed in Section 2.5 using (29) that, whenever $\Gamma > \Gamma_{\min}$, (10) is a necessary and sufficient condition for $C_{f,0}(\Gamma)$ to be positive, and we hence prove that if $C_{f,0}(\Gamma)$ is positive, then it is equal to (30).

To that end we first show that restricting X to be a function of U and S , i.e., $P_{U,X|S}$ to have the form (12), does not change the RHS of (30), nor does restricting the cardinality of \mathcal{U} to (13):

Lemma J.1. *Given a channel $W(y|x, s)$ and a PMF P_S on \mathcal{S} , consider*

$$\max_{\substack{P_{U,X|S}: \\ \mathbb{E}[\gamma(X)] \leq \Gamma}} \min_{\substack{P_{Y|U,X,S}: \\ P_{Y|U=u,X,S} \in \mathcal{P}(W), \forall u \in \mathcal{U}}} I(U; Y) - I(U; S), \quad (333)$$

where the maximization is over all chance variables U of finite support, the expectation is computed w.r.t. the joint PMF $P_S \times P_{U,X|S}$, and the mutual informations are computed w.r.t. the joint PMF $P_S \times P_{U,X|S} \times P_{Y|U,X,S}$. Restricting X to be a function of U and S , i.e., $P_{U,X|S}$ to have the form

$$P_{U,X|S}(u, x|s) = P_{U|S}(u|s) \mathbb{1}_{x=g(u,s)}, \quad (334)$$

does not change (333). Nor does requiring that U take values in a set \mathcal{U} whose cardinality $|\mathcal{U}|$ satisfies

$$|\mathcal{U}| \leq |\mathcal{X}|^{|\mathcal{S}|}. \quad (335)$$

Proof. The proof is essentially that of Lemma D.1 in Appendix D. We first show that restricting X to be a function of U and S does not change (333). In the proof of

Lemma D.1 it is shown that (200) is equal to (204), and the same line of argument implies here that (333) is equal to

$$\max_{P_V, h(\cdot), P_{U|S}: \mathbb{E}[\gamma(X)] \leq \Gamma} \min_{P_{Y|U,X,S}: P_{Y|U=u,X,S} \in \mathcal{P}(W), \forall u \in \mathcal{U}} I(U; Y) - I(U; S), \quad (336)$$

where the maximization is over all chance variables V of finite support \mathcal{V} , functions $h: \mathcal{U} \times \mathcal{V} \times \mathcal{S} \rightarrow \mathcal{X}$, and conditional PMFs over a finite set \mathcal{U} for which

$$\mathbb{E}[\gamma(X) \leq \Gamma], \quad (337)$$

where the expectation is computed w.r.t. the joint PMF $P_S \times P_V \times P_{U,X|V,S}$ and $P_{U,X|V,S}$ is defined in (205); and where the mutual informations are computed w.r.t. the joint PMF $P_S \times P_V \times P_{U,X|V,S} \times P_{Y|U,X,S}$. Unlike the proof of Lemma D.1, where we fix any PMF P_V on \mathcal{V} , any function $h: \mathcal{U} \times \mathcal{V} \times \mathcal{S} \rightarrow \mathcal{X}$, and any conditional PMF $P_{U|S}$, here we fix any P_V , any $h: \mathcal{U} \times \mathcal{V} \times \mathcal{S} \rightarrow \mathcal{X}$, and any $P_{U|S}$ for which (337) holds w.r.t. $P_S \times P_V \times P_{U,X|V,S}$, where $P_{U,X|V,S}$ is defined in (205). The line of argument leading to (214) in the proof of Lemma D.1 then implies that restricting X to be a function of U and S does no change (333). To show that restricting the cardinality of \mathcal{U} to (335) does not change (333), we fix any conditional PMF $P_{U,X|S}$ of the form (334) for which (337) holds w.r.t. $P_S \times P_{U,X|S}$. The line of argument leading to (232) in the proof of Lemma D.1 then implies that restricting the cardinality of \mathcal{U} to (335) does not change (333). \square

Direct Part of Theorem 2.20. From Lemma J.1 it follows that it suffices to establish the direct part of Theorem 2.20 for the case where the cardinality of \mathcal{U} is restricted to (13). The direct part is essentially that of Theorem 2.4 but with the following two modifications: 1) During the first B blocks we choose k -types $\{P_{U,X,S}^{(b)}\}_{b \in [1:B]}$ w.r.t. which

$$\mathbb{E}[\gamma(X)] \leq \Gamma. \quad (338)$$

This will guarantee that

$$\frac{1}{Bk} \sum_{i=1}^{Bk} \gamma(X_i) \leq \Gamma. \quad (339)$$

2) We pad Block $B+1$ with as many symbols from the set \mathcal{X}' as are needed to guarantee that

$$\frac{1}{k'} \sum_{i=Bk+1}^{Bk+k'} \gamma(X_i) \leq \Gamma, \quad (340)$$

where k' denotes the length of Block $B+1$.

By (339) and (340) the channel inputs' average cost satisfies the cost constraint (27). Padding Block $(B + 1)$ to guarantee (340) increases its length by a factor of at most

$$\tau \triangleq \left\lceil \frac{\gamma_{\max} - \gamma_{\min}}{\Gamma - \gamma_{\min}} \right\rceil. \quad (341)$$

Consequently, also with the padding, the last block does not affect the rate of the code.

To show that the coding scheme asymptotically achieves any rate smaller than the RHS of (30), we can argue essentially as in the proof of the direct part of Theorem 2.4. We will show that, for every rate R smaller than the RHS of (30) and every sufficiently-large blocklength n , our coding scheme can convey nR bits error-free in n channel uses. It follows from (107), (110), (338), and (341) that if the positive integers n , B , k and $\epsilon > 0$ are such that (111a) holds and

$$nR \leq Bk \left(\min_{P_S} \max_{\substack{P_{U,X|S}: \\ \mathbb{E}[\gamma(X)] \leq \Gamma}} \min_{\substack{P_{Y|U,X,S}: \\ P_{Y|U=u,X,S} \in \mathcal{P}(W), \forall u \in \mathcal{U}}} I(U; Y) - I(S; Y) - \delta(\epsilon, k) \right), \quad (342)$$

then our coding scheme can convey nR bits error-free in

$$Bk + \tau \lceil k \log |\mathcal{U}| + B \log(1 + k) |\mathcal{S}| \rceil n_{\text{bit}} \quad (343)$$

channel uses. It thus remains to exhibit positive integers B , k and some $\epsilon > 0$ such that, for every sufficiently-large blocklength n , (111a) and (342) hold and

$$Bk + \tau \lceil k \log |\mathcal{U}| + B \log(1 + k) |\mathcal{S}| \rceil n_{\text{bit}} \leq n. \quad (344)$$

As we argue next, when n is sufficiently large we can choose

$$B = \lfloor \sqrt{n} \rfloor - \tau \lceil \log |\mathcal{U}| + \log(1 + \sqrt{n}) |\mathcal{S}| \rceil n_{\text{bit}}, \quad (345a)$$

$$k = \lfloor \sqrt{n} \rfloor, \quad (345b)$$

and we can choose any $\epsilon > 0$ for which

$$R + \epsilon < \min_{P_S} \max_{\substack{P_{U,X|S}: \\ \mathbb{E}[\gamma(X)] \leq \Gamma}} \min_{\substack{P_{Y|U,X,S}: \\ P_{Y|U=u,X,S} \in \mathcal{P}(W), \forall u \in \mathcal{U}}} I(U; Y) - I(S; Y). \quad (346)$$

Note that, whenever n is sufficiently large, B is positive and (344) is satisfied. To see that also (111a) and (342) hold whenever n is sufficiently large, we first observe from (345b) that k tends to infinity as n tends to infinity. This implies that (111a) holds whenever n is sufficiently large, and that $\delta(\epsilon, k)$ (which is defined in (106), where $\gamma_k = \gamma_k(|\mathcal{U}|, |\mathcal{X}|, |\mathcal{S}|, |\mathcal{Y}|)$ converges to zero as k tends to infinity) converges to ϵ as n tends to infinity. We next observe that (345) implies that Bk/n converges to one as n tends to infinity. This, combined with the fact that $\delta(\epsilon, k)$ converges to ϵ as n tends to infinity and with (346), implies that (342) holds whenever n is sufficiently large. \square

Converse Part of Theorem 2.20. From Lemma J.1 it follows that it suffices to establish the converse part of Theorem 2.20 for the case where \mathcal{U} is any finite set. The converse is similar to that of Theorem 2.4. Fix a finite set \mathcal{M} , a blocklength n , and an (n, \mathcal{M}) zero-error code with n encoding mappings

$$f_i: \mathcal{M} \times \mathcal{S}^n \times \mathcal{Y}^{i-1} \rightarrow \mathcal{X}, \quad i \in [1 : n] \quad (347)$$

and $|\mathcal{M}|$ disjoint decoding sets $\mathcal{D}_m \subseteq \mathcal{Y}^n$, $m \in \mathcal{M}$, where the code is chosen so that, with probability one, the channel inputs X^n satisfy the cost constraint (27). We will show that, for some chance variable U of finite support \mathcal{U} , the rate $\frac{1}{n} \log |\mathcal{M}|$ of the code is upper-bounded by the RHS of (30).

Draw M uniformly over \mathcal{M} , and denote its distribution P_M . Since the code is a zero-error code, and since, with probability one, the channel inputs X^n satisfy the cost constraint (27), the following two hold:

$$\mathbb{P}[Y^n \in \mathcal{D}_M] = 1, \quad (348a)$$

$$\mathbb{P}[\gamma^{(n)}(X^n) \leq \Gamma] = 1, \quad (348b)$$

where \mathbb{P} is the distribution (118) of (M, S^n, X^n, Y^n) induced by P_M , the state distribution Q , the encoding mappings (347), and the channel law $W(y|x, s)$. As in the converse of Theorem 2.4, fix any PMF \tilde{P}_S on \mathcal{S} and any collection of n conditional PMFs $\{\tilde{P}_{Y_i|M, Y^{i-1}, S_{i+1}^n, X_i, S_i}\}_{i \in [1:n]}$ satisfying (120). These PMFs induce the PMF $\tilde{P}_{M, S^n, X^n, Y^n}$ of (121) on $\mathcal{M} \times \mathcal{S}^n \times \mathcal{X}^n \times \mathcal{Y}^n$. Since this PMF satisfies $\tilde{P}_{M, S^n, X^n, Y^n} \ll \mathbb{P}$, (348) implies

$$\tilde{P}_{M, S^n, X^n, Y^n}[Y^n \in \mathcal{D}_M] = 1, \quad (349a)$$

$$\tilde{P}_{M, S^n, X^n, Y^n}[\gamma^{(n)}(X^n) \leq \Gamma] = 1. \quad (349b)$$

Note that the latter (349b) implies that

$$\frac{1}{n} \sum_{i=1}^n \mathbb{E}[\gamma(X_i)] \leq \Gamma, \quad (350)$$

where the expectation in the i -th summand is computed w.r.t. the PMF \tilde{P}_{X_i} induced by $\tilde{P}_{M, S^n, X^n, Y^n}$.

The line of argument leading to (129) in the converse of Theorem 2.4 implies that every choice of \tilde{P}_S and $\{\tilde{P}_{Y_i|M, Y^{i-1}, S_{i+1}^n, X_i, S_i}\}_{i \in [1:n]}$ gives rise to an upper bound

$$\frac{1}{n} \log |\mathcal{M}| \leq \frac{1}{n} \sum_{i=1}^n [I(U_i; Y_i) - I(U_i; S_i)], \quad (351)$$

where the chance variables $\{U_i\}_{i \in [1:n]}$ are defined in (128), and the mutual informations in the i -th summand are computed w.r.t. the joint PMF $\tilde{P}_{U_i, X_i, S_i, Y_i}$ induced by $\tilde{P}_{M, S^n, X^n, Y^n}$. We next exhibit a PMF \tilde{P}_S and a collection of conditional PMFs $\{\tilde{P}_{Y_i|M, Y^{i-1}, S_{i+1}^n, X_i, S_i}\}_{i \in [1:n]}$ satisfying (120) for which each summand on the RHS of (351) is upper-bounded by the RHS of (30).

We begin with the choice of $\{\tilde{P}_{Y_i|M, Y^{i-1}, S_{i+1}^n, X_i, S_i}\}_{i \in [1:n]}$. As in the converse of Theorem 2.4, choosing a collection of conditional PMFs $\{\tilde{P}_{Y_i|M, Y^{i-1}, S_{i+1}^n, X_i, S_i}\}_{i \in [1:n]}$ that satisfy (120) is tantamount to choosing a collection of conditional PMFs $\{\tilde{P}_{Y_i|U_i, X_i, S_i}\}_{i \in [1:n]}$ that satisfy (131). We shall choose the latter collection, and we shall do so as in the converse of Theorem 2.4. Consequently, our choice of $\tilde{P}_{Y_i|U_i, X_i, S_i}$ can depend not only on our choice of \tilde{P}_S and our previous choices of $\{\tilde{P}_{Y_j|U_j, X_j, S_j}\}_{j \in [1:i-1]}$ but also on $\tilde{P}_{U_i, X_i, S_i}$, and hence we can choose $\tilde{P}_{Y_i|U_i, X_i, S_i}$ as one that—among all conditional PMFs satisfying (131)—minimizes

$$I(U_i; Y_i) - I(U_i; S_i), \quad (352)$$

where the mutual informations are computed w.r.t. the joint PMF $\tilde{P}_{U_i, X_i, S_i} \times \tilde{P}_{Y_i|U_i, X_i, S_i}$. Because (121) implies that

$$\tilde{P}_{S_i} = \tilde{P}_S, \quad i \in [1 : n] \quad (353)$$

and by (350), which holds because the chosen conditional PMFs $\{\tilde{P}_{Y_i|U_i, X_i, S_i}\}_{i \in [1:n]}$ satisfy (131), we find that, for our choice of $\{\tilde{P}_{Y_i|U_i, X_i, S_i}\}_{i \in [1:n]}$,

$$\begin{aligned} & \frac{1}{n} \sum_{i=1}^n \left[I(U_i; Y_i) - I(U_i; S_i) \right] \\ & \leq \max_{\{\tilde{P}_{U_i, X_i|S_i}\}_{i \in [1:n]} : \frac{1}{n} \sum_{i=1}^n \mathbb{E}[\gamma(X_i)] \leq \Gamma} \min_{\{\tilde{P}_{Y_i|U_i, X_i, S_i}\}_{i \in [1:n]} : \tilde{P}_{Y_i|U_i=u_i, X_i, S_i} \in \mathcal{P}(W), \forall u_i \in \mathcal{U}_i} \frac{1}{n} \sum_{i=1}^n \left[I(U_i; Y_i) - I(U_i; S_i) \right], \end{aligned} \quad (354)$$

where the expectation in the i -th summand is computed w.r.t. the joint PMF $\tilde{P}_{S_i} \times \tilde{P}_{U_i, X_i|S_i}$, and the mutual informations in the i -th summand are computed w.r.t. the joint PMF $\tilde{P}_{S_i} \times \tilde{P}_{U_i, X_i|S_i} \times \tilde{P}_{Y_i|U_i, X_i, S_i}$. The chosen conditional PMFs $\{\tilde{P}_{Y_i|U_i, X_i, S_i}\}_{i \in [1:n]}$ satisfy (131), and hence (351) and (354) imply that

$$\begin{aligned} & \frac{1}{n} \log |\mathcal{M}| \\ & \leq \max_{\{\tilde{P}_{U_i, X_i|S_i}\}_{i \in [1:n]} : \frac{1}{n} \sum_{i=1}^n \mathbb{E}[\gamma(X_i)] \leq \Gamma} \min_{\{\tilde{P}_{Y_i|U_i, X_i, S_i}\}_{i \in [1:n]} : \tilde{P}_{Y_i|U_i=u_i, X_i, S_i} \in \mathcal{P}(W), \forall u_i \in \mathcal{U}_i} \frac{1}{n} \sum_{i=1}^n \left[I(U_i; Y_i) - I(U_i; S_i) \right], \end{aligned} \quad (355)$$

where the expectation in the i -th summand is computed w.r.t. the joint PMF $\tilde{P}_{S_i} \times \tilde{P}_{U_i, X_i | S_i}$, and the mutual informations in the i -th summand are computed w.r.t. the joint PMF $\tilde{P}_{S_i} \times \tilde{P}_{U_i, X_i | S_i} \times \tilde{P}_{Y_i | U_i, X_i, S_i}$.

By the definition of U_i (128) the cardinality of the support \mathcal{U}_i of U_i satisfies (138). Consequently, (353) and (355) imply that

$$\begin{aligned} & \frac{1}{n} \log |\mathcal{M}| \\ & \leq \max_{\substack{\tilde{P}_{U, X | V, S}: \\ \mathbb{E}[\gamma(X)] \leq \Gamma}} \min_{\substack{\tilde{P}_{Y | V, U, X, S}: \\ \tilde{P}_{Y | (V, U) = (i, u), X, S} \in \mathcal{P}(W), \forall (i, u) \in [1:n] \times \mathcal{U}}} I(U; Y | V) - I(U; S | V) \end{aligned} \quad (356)$$

$$\begin{aligned} & \leq \max_{\substack{\tilde{P}_{U, X | V, S}: \\ \mathbb{E}[\gamma(X)] \leq \Gamma}} \min_{\substack{\tilde{P}_{Y | V, U, X, S}: \\ \tilde{P}_{Y | (V, U) = (i, u), X, S} \in \mathcal{P}(W), \forall (i, u) \in [1:n] \times \mathcal{U}}} I(V, U; Y) - I(V, U; S), \end{aligned} \quad (357)$$

where V is a time-sharing random-variable that is drawn uniformly over $[1 : n]$ and U an auxiliary chance variable taking values in a finite set \mathcal{U} ; where the mutual informations are computed w.r.t. the joint PMF $\tilde{P}_S \times P_V \times \tilde{P}_{U, X | V, S} \times \tilde{P}_{Y | V, U, X, S}$; and where the second inequality holds because conditioning cannot increase entropy, and because S and V are independent under $\tilde{P}_S \times P_V \times \tilde{P}_{U, X | V, S} \times \tilde{P}_{Y | V, U, X, S}$. By defining the auxiliary chance variable $\tilde{U} = (U, V)$, we obtain from (357) that every choice of \tilde{P}_S gives rise to an upper bound

$$\begin{aligned} & \frac{1}{n} \log |\mathcal{M}| \\ & \leq \max_{\substack{\tilde{P}_{\tilde{U}, X | S}: \\ \mathbb{E}[\gamma(X)] \leq \Gamma}} \min_{\substack{\tilde{P}_{Y | \tilde{U}, X, S}: \\ \tilde{P}_{Y | \tilde{U} = \tilde{u}, X, S} \in \mathcal{P}(W), \forall \tilde{u} \in \tilde{\mathcal{U}}}} I(\tilde{U}; Y) - I(\tilde{U}; S), \end{aligned} \quad (358)$$

where \tilde{U} is an auxiliary chance variable taking values in a finite set $\tilde{\mathcal{U}}$, and the mutual informations are computed w.r.t. the joint PMF $\tilde{P}_S \times \tilde{P}_{\tilde{U}, X | S} \times \tilde{P}_{Y | \tilde{U}, X, S}$.

Having established (358), we are now ready to conclude the proof of the converse. Since we can choose any PMF \tilde{P}_S on \mathcal{S} , we can choose one that—among all PMFs on \mathcal{S} —yields the tightest bound, i.e., minimizes

$$\begin{aligned} & \max_{\substack{\tilde{P}_{\tilde{U}, X | S}: \\ \mathbb{E}[\gamma(X)] \leq \Gamma}} \min_{\substack{\tilde{P}_{Y | \tilde{U}, X, S}: \\ \tilde{P}_{Y | \tilde{U} = \tilde{u}, X, S} \in \mathcal{P}(W), \forall \tilde{u} \in \tilde{\mathcal{U}}}} I(\tilde{U}; Y) - I(\tilde{U}; S), \end{aligned} \quad (359)$$

where \tilde{U} is an auxiliary chance variable taking values in a finite set $\tilde{\mathcal{U}}$, and the mutual informations are computed w.r.t. the joint PMF $\tilde{P}_S \times \tilde{P}_{\tilde{U}, X | S} \times \tilde{P}_{Y | \tilde{U}, X, S}$. For this choice

of \tilde{P}_S (358) implies that

$$\begin{aligned} & \frac{1}{n} \log |\mathcal{M}| \\ & \leq \min_{\tilde{P}_S} \max_{\substack{\tilde{P}_{\tilde{U},X|S}: \\ \mathbb{E}[\gamma(X)] \leq \Gamma}} \min_{\substack{\tilde{P}_{Y|\tilde{U},X,S}: \\ \tilde{P}_{Y|\tilde{U}=\tilde{u},X,S} \in \mathcal{P}(W), \forall \tilde{u} \in \tilde{\mathcal{U}}}} I(\tilde{U}; Y) - I(\tilde{U}; S), \end{aligned} \quad (360)$$

where \tilde{U} is an auxiliary chance variable taking values in a finite set $\tilde{\mathcal{U}}$, and the mutual informations are computed w.r.t. the joint PMF $\tilde{P}_S \times \tilde{P}_{\tilde{U},X|S} \times \tilde{P}_{Y|\tilde{U},X,S}$. \square

K A Proof of Remark 2.22

Proof. Fix some PMF P_X on \mathcal{X} , and define the function

$$\begin{aligned} \rho: \mathcal{P}(W) &\rightarrow \mathbb{R}_0^+ \\ V &\mapsto I(P_X, V). \end{aligned} \quad (361)$$

To prove (34), we will show that every V that minimizes $\rho(\cdot)$ satisfies

$$\rho(V) \geq \min_{y \in \mathcal{Y}} -\log \sum_{x \in \mathcal{X}: W(y|x) > 0} P_X(x). \quad (362)$$

From this we will then obtain (34) by maximizing both sides over all choices of P_X for which $\mathbb{E}[\gamma(X)] \leq \Gamma$. This will conclude the proof of Remark 2.22, because Example 2.23 demonstrates that Inequality (34) can be strict.

To show that every minimizer of $\rho(\cdot)$ satisfies (362), we first establish that $V \in \mathcal{P}(W)$ minimizes $\rho(\cdot)$ only if

$$\frac{V(y|x)}{(P_X V)(y)} = \frac{V(y'|x)}{(P_X V)(y')}, \quad \left(\forall (x, y, y') \in \mathcal{X} \times \mathcal{Y} \times \mathcal{Y} \text{ s.t. } V(y|x) V(y'|x) > 0 \right). \quad (363)$$

We prove the contrapositive: we show that if for some $V \in \mathcal{P}(W)$

$$\begin{aligned} & \exists (x, y, y') \in \mathcal{X} \times \mathcal{Y} \times \mathcal{Y} \text{ s.t.} \\ & V(y|x) V(y'|x) > 0 \quad \text{and} \quad \frac{V(y|x)}{(P_X V)(y)} < \frac{V(y'|x)}{(P_X V)(y')}, \end{aligned} \quad (364)$$

then V cannot be a minimizer of $\rho(\cdot)$. Our proof is by contradiction. To reach a contradiction, suppose that $V \in \mathcal{P}(W)$ minimizes $\rho(\cdot)$ and (364) holds. Since

$$I(P_X, V) = \sum_{(x,y) \in \mathcal{X} \times \mathcal{Y}} P_X(x) V(y|x) \log \frac{V(y|x)}{(P_X V)(y)}, \quad (365)$$

it follows that for every $(x, y) \in \mathcal{X} \times \mathcal{Y}$

$$\frac{\partial \rho}{\partial V(y|x)} = P_X(x) \log \frac{V(y|x)}{(P_X V)(y)}. \quad (366)$$

This and (364) imply that for all sufficiently-small δ , and a fortiori for some δ satisfying

$$0 < \delta \leq (1 - V(y|x)) \wedge V(y'|x), \quad (367)$$

$\rho(\cdot)$ decreases when we replace $V(y|x)$ by $V(y|x) + \delta$ and $V(y'|x)$ by $V(y'|x) - \delta$. This contradicts our assumption that V minimizes $\rho(\cdot)$, because replacing $V(y|x)$ by $V(y|x) + \delta$ and $V(y'|x)$ by $V(y'|x) - \delta$ yields some transition law V' in $\mathcal{P}(W)$. (The transition law V' is in $\mathcal{P}(W)$, because $V \in \mathcal{P}(W)$ and by (367).) This contradiction proves that (363) is a necessary condition for $V \in \mathcal{P}(W)$ to minimize $\rho(\cdot)$.

Having proved the necessity of (363), we are now ready to establish (362). To that end let

$$\gamma = \max_{y \in \mathcal{Y}} \sum_{x \in \mathcal{X}: W(y|x) > 0} P_X(x), \quad (368)$$

and fix some transition law $V \in \mathcal{P}(W)$ that minimizes $\rho(\cdot)$ and for which (363) hence holds. By (363) there exist $\{\alpha_x\}_{x \in \mathcal{X}}$ satisfying that, whenever $V(y|x) > 0$,

$$\alpha_x = \frac{V(y|x)}{(P_X V)(y)}. \quad (369)$$

Consequently,

$$\rho(V) = \sum_{x \in \mathcal{X}} P_X(x) \sum_{y \in \mathcal{Y}: V(y|x) > 0} (P_X V)(y) \alpha_x \log \alpha_x \quad (370)$$

$$= \sum_{y \in \mathcal{Y}} (P_X V)(y) \sum_{x \in \mathcal{X}: V(y|x) > 0} P_X(x) \alpha_x \log \alpha_x \quad (371)$$

$$\stackrel{(a)}{\geq} \sum_{y \in \mathcal{Y}} (P_X V)(y) \left(\sum_{x'' \in \mathcal{X}: V(y|x'') > 0} P_X(x'') \right) \times \frac{\sum_{x \in \mathcal{X}: V(y|x) > 0} P_X(x) \alpha_x}{\sum_{x'' \in \mathcal{X}: V(y|x'') > 0} P_X(x'')} \log \frac{\sum_{x \in \mathcal{X}: V(y|x) > 0} P_X(x) \alpha_x}{\sum_{x'' \in \mathcal{X}: V(y|x'') > 0} P_X(x'')} \quad (372)$$

$$\stackrel{(b)}{=} \sum_{y \in \mathcal{Y}} (P_X V)(y) \log \frac{1}{\sum_{x'' \in \mathcal{X}: V(y|x'') > 0} P_X(x'')} \quad (373)$$

$$\stackrel{(c)}{\geq} - \sum_{y \in \mathcal{Y}} (P_X V)(y) \log \gamma \quad (374)$$

$$= - \log \gamma \quad (375)$$

$$\stackrel{(d)}{=} \min_{y \in \mathcal{Y}} - \log \sum_{x \in \mathcal{X}: W(y|x) > 0} P_X(x), \quad (376)$$

where (a) holds because the function

$$\xi \mapsto \xi \log \xi, \quad \xi \in \mathbb{R}^+$$

is convex; (b) holds because (369) holds whenever $V(y|x) > 0$; (c) holds because $V \in \mathcal{P}(W)$ and (368) combine to imply that

$$\sum_{x \in \mathcal{X}: V(y|x) > 0} P_X(x) \leq \sum_{x \in \mathcal{X}: W(y|x) > 0} P_X(x) \leq \gamma, \quad y \in \mathcal{Y}; \quad (377)$$

and (d) holds by (368). Inequality (376) concludes the proof of (362). \square

L Analysis of Example 2.25

For the SD-DMC $W(y|x, s)$ of Example 2.25 we show that, subject to the cost constraint (43) with $\Gamma > 0$ satisfying (46), the zero-error capacity with acausal SI is positive. Given some blocklength n , some message set \mathcal{M} , and some encoding mapping

$$f: \mathcal{M} \times \mathcal{S}^n \rightarrow \mathcal{X}^n, \quad (378)$$

let $\mathbf{y}(m, \mathbf{s})$ denote the output sequence that is produced when the transmitter uses the encoding mapping (378) to convey Message m and the channel-state sequence is \mathbf{s} . Let $\mathcal{S}^n(\Lambda)$ denote the set of n -length state-sequences of highest allowed cost

$$\mathcal{S}^n(\Lambda) = \{\mathbf{s} \in \mathcal{S}^n: n\lambda^{(n)}(\mathbf{s}) = \lfloor \Lambda n \rfloor\}. \quad (379)$$

We begin with the following two observations: 1) From Table 2 we see that, if

$$\{\tilde{\mathbf{y}}(m, \mathbf{s})\}_{(m, \mathbf{s}) \in \mathcal{M} \times \mathcal{S}^n} \subseteq \mathcal{Y}^n \quad (380)$$

is such that for every $\mathbf{s} \in \mathcal{S}^n$

$$\left((s_i = 1) \implies (\tilde{y}_i(m, \mathbf{s}) = 1) \right), \quad \forall (m, i) \in \mathcal{M} \times [1 : n], \quad (381)$$

then there exists an encoding mapping f of the form (378) for which

$$\mathbf{y}(m, \mathbf{s}) = \tilde{\mathbf{y}}(m, \mathbf{s}), \quad \forall (m, \mathbf{s}) \in \mathcal{M} \times \mathcal{S}^n. \quad (382)$$

2) From the definition of $\mathcal{S}^n(\Lambda)$ it follows that, if $\mathbf{s} \in \mathcal{S}^n$ is such that $n\lambda^{(n)}(\mathbf{s}) < \lfloor \Lambda n \rfloor$, then there exists some $\mathbf{s}' \in \mathcal{S}^n(\Lambda)$ satisfying

$$\left((s_i = 1) \implies (s'_i = 1) \right), \quad \forall i \in [1 : n]. \quad (383)$$

For such \mathbf{s}' , any binary n -tuple $\tilde{\mathbf{y}}(m, \mathbf{s}')$ satisfying

$$\left((s'_i = 1) \implies \left(\tilde{y}_i(m, \mathbf{s}') = 1 \right) \right), \quad \forall i \in [1 : n] \quad (384)$$

also satisfies

$$\left((s_i = 1) \implies \left(\tilde{y}_i(m, \mathbf{s}') = 1 \right) \right), \quad \forall i \in [1 : n]. \quad (385)$$

These two observations imply that to every collection

$$\{\tilde{\mathbf{y}}(m, \mathbf{s})\}_{(m, \mathbf{s}) \in \mathcal{M} \times \mathcal{S}^n(\Lambda)} \subseteq \mathcal{Y}^n \quad (386)$$

that satisfies (381) for every $\mathbf{s} \in \mathcal{S}^n(\Lambda)$ there corresponds an encoding mapping of the form (378) for which: 1) for every $\mathbf{s} \in \mathcal{S}^n(\Lambda)$

$$\mathbf{y}(m, \mathbf{s}) = \tilde{\mathbf{y}}(m, \mathbf{s}), \quad \forall m \in \mathcal{M}; \quad (387a)$$

and 2) for every $\mathbf{s} \in \mathcal{S}^n$ for which $n\lambda^{(n)}(\mathbf{s}) < \lfloor \Lambda n \rfloor$

$$\exists \mathbf{s}' \in \mathcal{S}^n(\Lambda) \text{ s.t. } \left(\mathbf{y}(m, \mathbf{s}) = \tilde{\mathbf{y}}(m, \mathbf{s}'), \quad \forall m \in \mathcal{M} \right). \quad (387b)$$

The state sequence $\mathbf{s} \in \mathcal{S}^n$ satisfies the cost constraint (43) if $n\lambda^{(n)}(\mathbf{s}) \leq \lfloor \Lambda n \rfloor$. Consequently, if the collection in (386)—in addition to satisfying (381) for every $\mathbf{s} \in \mathcal{S}^n(\Lambda)$ —also satisfies that

$$\left(\left((m, \mathbf{s}) \neq (m', \mathbf{s}') \right) \implies \left(\tilde{\mathbf{y}}(m, \mathbf{s}) \neq \tilde{\mathbf{y}}(m', \mathbf{s}') \right) \right), \quad \forall (m, \mathbf{s}), (m', \mathbf{s}') \in \mathcal{M} \times \mathcal{S}^n(\Lambda), \quad (388)$$

then we obtain from (387) that the encoding mapping f corresponding to the collection and the decoding sets

$$\mathcal{D}_m = \bigcup_{\mathbf{s} \in \mathcal{S}^n(\Lambda)} \{\tilde{\mathbf{y}}(m, \mathbf{s})\}, \quad m \in \mathcal{M} \quad (389)$$

constitute an (n, \mathcal{M}) zero-error code for our channel under the cost constraint (43).

To show that under the cost constraint (43) the zero-error capacity with acausal SI is positive, it thus suffices to exhibit some positive rate $R > 0$ for which for every sufficiently-large n there exists some finite set \mathcal{M} of cardinality $|\mathcal{M}| \geq 2^{nR}$ and some collection of $|\mathcal{M}| |\mathcal{S}^n(\Lambda)|$ distinct binary n -tuples (386) that satisfies (381) for every $\mathbf{s} \in \mathcal{S}^n(\Lambda)$.

To that end we first note that the cardinality of $\mathcal{S}^n(\Lambda)$ is upper-bounded by

$$|\mathcal{S}^n(\Lambda)| = \binom{n}{\lfloor \Lambda n \rfloor} \leq 2^{n h_b(\Lambda)}, \quad (390)$$

where we used the inequality $h_b(\lfloor \Lambda n \rfloor / n) \leq h_b(\Lambda)$ (which holds because $\Lambda < 1/2$). We also note that for every state sequence $\mathbf{s} \in \mathcal{S}^n(\Lambda)$ there exist $2^{\lceil n(1-\Lambda) \rceil}$ binary n -tuples $\tilde{\mathbf{y}}$ that satisfy

$$\left((s_i = 1) \implies (\tilde{y}_i = 1) \right), \quad \forall i \in [1 : n]. \quad (391)$$

We now construct a collection of $|\mathcal{M}| |\mathcal{S}^n(\Lambda)|$ distinct binary n -tuples (386) that satisfies (381) for every $\mathbf{s} \in \mathcal{S}^n(\Lambda)$ as follows. We sequentially allocate to each pair $(m, \mathbf{s}) \in \mathcal{M} \times \mathcal{S}^n(\Lambda)$ some $\tilde{\mathbf{y}}(m, \mathbf{s})$ from the binary n -tuples $\tilde{\mathbf{y}}$ that satisfy (391) and that have not yet been allocated to some other pair (m', \mathbf{s}') . There are $|\mathcal{M}| |\mathcal{S}^n(\Lambda)| - 1$ such other pairs (m', \mathbf{s}') to which we may or may not have allocated some $\tilde{\mathbf{y}}(m', \mathbf{s}')$ yet, and there are at least $2^{\lceil n(1-\Lambda) \rceil}$ binary n -tuples $\tilde{\mathbf{y}}$ that satisfy (391). Consequently, at most $|\mathcal{M}| |\mathcal{S}^n(\Lambda)| - 1$ binary n -tuples could have already been allocated, and if

$$|\mathcal{M}| |\mathcal{S}^n(\Lambda)| - 1 < 2^{\lceil n(1-\Lambda) \rceil}, \quad (392)$$

then there is at least one binary n -tuples $\tilde{\mathbf{y}}$ that satisfies (391) and that has not been allocated yet. Hence, if (392) holds, then our construction produces a collection of $|\mathcal{M}| |\mathcal{S}^n(\Lambda)|$ distinct binary n -tuples (386) that satisfies (381) for every $\mathbf{s} \in \mathcal{S}^n(\Lambda)$. From (390) we obtain that (392) holds whenever

$$|\mathcal{M}| \leq 2^{n(1-\Lambda-h_b(\Lambda))}, \quad (393)$$

and hence every positive rate $R > 0$ satisfying

$$R \leq 1 - \Lambda - h_b(\Lambda) \quad (394)$$

is achievable. This, combined with (46), implies that under the cost constraint (43) the zero-error capacity with acausal SI is positive.

M A Proof of Theorem 2.26

Lemma D.1 in Appendix D implies that restricting X to be a function of U and S , i.e., $P_{U,X|S}$ to have the form (12), does not change the RHS of (50), nor does restricting the cardinality of \mathcal{U} to (13). To prove Theorem 2.26 it thus suffices to establish a direct part for the case where \mathcal{U} is restricted to (13) and a converse part for the case where \mathcal{U} is any finite set. We first establish the direct part.

Direct Part. We assume that (49) holds and show that the RHS of (50) is achievable. The necessity of (49) is part of the converse. If the RHS of (50) is zero, then there is nothing to prove, so we assume that it is positive. To prove that the RHS of (50) is achievable, we shall show that for every $l \in \mathbb{N}$ the RHS of (50) is a lower bound for $C_{f,0}^{(2)}(\Lambda, l)$. To that end fix any $l \in \mathbb{N}$. The proof builds on the proofs of Remark 3.1 and the direct part of Theorem 2.4, adapting both to the state constraint (48). We partition the blocklength- n transmission into $B + 2$ blocks, with each of the first B blocks being of length k , where k is a multiple of l ; with Block $(B + 1)$ being of length k' ; and with Block $(B + 2)$ being of length $n - Bk - k'$. The only purpose of Block $(B + 2)$ is to allow $Bk + k'$ to be smaller than n : in this block the encoder can thus transmit arbitrary inputs with the decoder ignoring the corresponding outputs. The choice we shall later make for k and k' will be such that the last two blocks be of negligible length compared to Bk and therefore not affect the code's asymptotic rate.

Before the transmission begins, the encoder is revealed the realization $\mathbf{s}^{(b)} \triangleq S_{(b-1)k+1}^{bk}$ of the Block- b state-sequence for every $b \in [1 : B]$ and the realization $\mathbf{s}^{(B+1)} \triangleq S_{Bk+1}^{Bk+k'}$ of the Block- $(B + 1)$ state-sequence. In the first B blocks our scheme draws on the scheme we used in the direct part of Theorem 2.4. But instead of reducing the set of messages of positive posterior probability given the channel outputs, in the present setting we consider pairs of messages and possible Block- $(B + 1)$ state-sequences, and each of the blocks 1 through B reduces the set of such pairs that have a positive posterior probability given the channel outputs. For every $b \in [1 : B]$ we thus adapt the Block b transmission as follows. Because k is a multiple of l , the cost constraint (48) implies that in the first B blocks

$$\sum_{s \in \mathcal{S}} P_{\mathbf{s}^{(b)}}(s) \lambda(S) \leq \Lambda, \quad \forall b \in [1 : B], \quad (395a)$$

and in Block $(B + 1)$

$$\frac{1}{l} \sum_{i=(j-1)l+1}^{jl} \lambda(s_i^{(B+1)}) \leq \Lambda, \quad (\forall j \in \mathbb{N} \text{ s.t. } jl \leq k'). \quad (395b)$$

Assume for now that the decoder—while incognizant of $\mathbf{s}^{(1)}, \dots, \mathbf{s}^{(B)}$ —knows the empirical types $P_{\mathbf{s}^{(1)}}, \dots, P_{\mathbf{s}^{(B)}}: \text{Block } (B + 1)$ will ensure that the scheme works even though the decoder is incognizant of these types. Let $\mathcal{I}_0 \subseteq \mathcal{M} \times \mathcal{S}^{k'}$ be the set of all possible pairs of message $m' \in \mathcal{M}$ and Block- $(B + 1)$ state-sequence $\mathbf{s}' \in \mathcal{S}^{k'}$ satisfying (395b), i.e.,

$$\frac{1}{l} \sum_{i=(j-1)l+1}^{jl} \lambda(s'_i) \leq \Lambda, \quad (\forall j \in \mathbb{N} \text{ s.t. } jl \leq k'), \quad (396)$$

and let \mathcal{I}_b be the post-Block- b ambiguity-set, i.e., the (random) subset of \mathcal{I}_{b-1} comprising the elements in \mathcal{I}_{b-1} of positive posterior probability given the Block- b outputs $\mathbf{y}^{(b)} \triangleq Y_{(b-1)k+1}^{bk}$ and the empirical type $P_{\mathbf{s}^{(b)}}$. Choose some k -type $P_{U,X,S}^{(b)}$ whose \mathcal{S} -marginal $P_S^{(b)}$ equals $P_{\mathbf{s}^{(b)}}$, fix some $\epsilon > 0$, and define Θ as in (91). In the following, unless otherwise specified, all entropies and mutual informations are computed w.r.t. the joint PMF $P_{U,X,S}^{(b)}$. Unlike the scheme we used in the direct part of Theorem 2.4, where it was the survivor set \mathcal{M}_{b-1} that was partitioned into Θ subsets, here it is the ambiguity set \mathcal{I}_{b-1} that is partitioned into Θ subsets. The arguments leading to (99) in the direct part of Theorem 2.4 then imply that we can find a positive integer $\eta_0 = \eta_0(|\mathcal{X}|, |\mathcal{S}|, \epsilon)$ that guarantees that, for every $k \geq \eta_0$,

$$|\mathcal{I}_b| \leq \left(\max_{\substack{P_{Y|U,X,S}: \\ P_{Y|U=u,X,S} \in \mathcal{P}(W), \forall u \in \mathcal{U}}} 2^{-k(I(U;Y) - I(U;S) - (\epsilon + \beta_k))} \right) |\mathcal{I}_{b-1}|, \quad (397a)$$

whenever

$$|\mathcal{I}_{b-1}| \geq 2^{k \log |\mathcal{U}|}, \quad (397b)$$

where the mutual informations are computed w.r.t. the joint PMF $P_{U,X,S}^{(b)} \times P_{Y|U,X,S}$, and where β_k is defined in (94) and hence converges to zero as k tends to infinity.

Since we can choose any k -type $P_{U,X,S}^{(b)}$ whose \mathcal{S} -marginal $P_S^{(b)}$ is $P_{\mathbf{s}^{(b)}}$, we can choose $P_{U,X,S}^{(b)} = P_{\mathbf{s}^{(b)}} \times P_{U,X|S}^{(b)}$, where $P_{U,X|S}^{(b)}$ is the conditional k -type that—among all conditional k -types—maximizes

$$\min_{\substack{P_{Y|U,X,S}: \\ P_{Y|U=u,X,S} \in \mathcal{P}(W), \forall u \in \mathcal{U}}} I(U;Y) - I(U;S), \quad (398)$$

where the mutual informations are computed w.r.t. the joint PMF $P_{\mathbf{s}^{(b)}} \times P_{U,X|S}^{(b)} \times P_{Y|U,X,S}$. Every conditional PMF can be approximated in the total variation distance by a conditional k -type when k is sufficiently large; and, because entropy and mutual information are continuous in this distance [8, Lemma 2.7], it follows that—for the above choice of the conditional k -type and some $\gamma_k = \gamma_k(|\mathcal{U}|, |\mathcal{X}|, |\mathcal{S}|, |\mathcal{Y}|)$, which converges to zero as k tends to infinity—(395a) and (397) imply that when $|\mathcal{I}_{b-1}| \geq 2^{k \log |\mathcal{U}|}$

$$|\mathcal{I}_b| \leq \left(\max_{\substack{P_S: \\ \mathbb{E}[\lambda(S)] \leq \Lambda}} \min_{P_{U,X|S}} \max_{\substack{P_{Y|U,X,S}: \\ P_{Y|U=u,X,S} \in \mathcal{P}(W), \forall u \in \mathcal{U}}} 2^{-k(I(U;Y) - I(U;S) - \epsilon + \gamma_k)} \right) |\mathcal{I}_{b-1}|, \quad (399)$$

where the mutual informations are computed w.r.t. the joint PMF $P_S \times P_{U,X|S} \times P_{Y|U,X,S}$. Because our scheme works for any $\epsilon > 0$, it follows that for every $\epsilon > 0$ and positive integer $k \geq \eta_0(|\mathcal{X}|, |\mathcal{S}|, \epsilon)$ each of Blocks 1 through B is guaranteed to reduce the ambiguity

set by a factor of at least

$$\max_{\substack{P_S: \\ \mathbb{E}[\lambda(S)] \leq \Lambda}} \min_{P_{U,X|S}} \max_{\substack{P_{Y|U,X,S}: \\ P_{Y|U=u,X,S} \in \mathcal{P}(W), \forall u \in \mathcal{U}}} 2^{-k(I(U;Y) - I(U;S) - \delta(\epsilon, k))}, \quad (400)$$

until $|\mathcal{I}_B|$ is smaller than $2^{k \log |\mathcal{U}|}$. Here the mutual informations are computed w.r.t. the joint PMF $P_S \times P_{U,X|S} \times P_{Y|U,X,S}$, and $\delta(\epsilon, k)$ is defined in (106) and hence converges to zero as ϵ tends to zero and k to infinity.

Since we assume that the RHS of (50) is positive; and, because $\delta(\epsilon, k)$ converges to zero as $\epsilon \downarrow 0$ and $k \rightarrow \infty$, it follows that we can choose ϵ sufficiently small and B and k sufficiently large so that

$$k \geq \eta_0(|\mathcal{X}|, |\mathcal{S}|, \epsilon) \quad (401a)$$

and

$$\left(\max_{\substack{P_S: \\ \mathbb{E}[\lambda(S)] \leq \Lambda}} \min_{P_{U,X|S}} \max_{\substack{P_{Y|U,X,S}: \\ P_{Y|U=u,X,S} \in \mathcal{P}(W), \forall u \in \mathcal{U}}} 2^{-Bk(I(U;Y) - I(U;S) - \delta(\epsilon, k))} \right) |\mathcal{M}| |\mathcal{S}|^{k'} \leq 2^{k \log |\mathcal{U}|}. \quad (401b)$$

This guarantees that

$$|\mathcal{I}_B| \leq 2^{k \log |\mathcal{U}|}, \quad (402)$$

because each block reduces the ambiguity set by the factor in (400) until $|\mathcal{I}_B|$ is smaller than $2^{k \log |\mathcal{U}|}$.

We now deal with Block $(B+1)$. Because the decoder is incognizant of the empirical types $\{P_{\mathbf{s}^{(b)}}\}_{b \in [1:B]}$, it cannot compute the post-Block- B ambiguity-set \mathcal{I}_B . The uncertainty that needs to be addressed is about the message, the Block- $(B+1)$ state-sequence, as well as the B empirical types of $\mathbf{s}^{(1)}, \dots, \mathbf{s}^{(B)}$. Let $\mathcal{J}_B \subseteq \mathcal{M} \times \mathcal{S}^{k'}$ denote the union of the post-Block- B ambiguity-sets corresponding to all the different B -tuples of k -types on \mathcal{S} , i.e., \mathcal{J}_B is the set of pairs of messages and possible Block- $(B+1)$ state-sequences that have a positive posterior probability given only the outputs $\{\mathbf{y}^{(b)}\}_{b \in [1:B]}$ (and not the k -types $\{P_{\mathbf{s}^{(b)}}\}_{b \in [1:B]}$). Because the post-Block- B ambiguity-set corresponding to any given B -tuple of k -types on \mathcal{S} satisfies (402), and because there are at most $(1+k)^B |\mathcal{S}|$ B -tuples of k -types on \mathcal{S}

$$|\mathcal{J}_B| \leq 2^{k \log |\mathcal{U}| + B \log(1+k) |\mathcal{S}|}. \quad (403)$$

In Block $(B+1)$ we resolve the set \mathcal{J}_B . This will guarantee that the decoder can recover the transmitted message m error-free.

Block $(B + 1)$ is similar to Phase 2 of the scheme we used to prove Remark 3.1: the encoder allocates to every pair $(m', s') \in \mathcal{J}_B$ a length- k' codeword $\mathbf{x}(m', s')$, where the codewords are chosen so that

$$\begin{aligned} & \left(\forall (m', s'), (m'', s'') \in \mathcal{J}_B \text{ s.t. } m' \neq m'' \right) \quad \exists i \in [1 : k'] \text{ s.t.} \\ & \left(W(y|x_i(m', s'), s'_i) W(y|x_i(m'', s''), s''_i) = 0, \forall y \in \mathcal{Y} \right). \end{aligned} \quad (404)$$

(We shall shortly use a random coding argument to show that this can be done.) To convey the message m , the encoder transmits in Block $(B + 1)$ the codeword $\mathbf{x}(m, \mathbf{s}^{(B+1)})$. Condition (404) implies that, upon observing the Block- $(B + 1)$ outputs $\mathbf{y}^{(B+1)} \triangleq Y_{Bk+1}^{Bk+k'}$, the decoder, who knows \mathcal{J}_B and the codewords $\{\mathbf{x}(m', s')\}$, can determine the transmitted message m error-free, because, for the true realization $\mathbf{s}^{(B+1)}$ of the Block- $(B + 1)$ state-sequence,

$$\prod_{i=1}^{k'} W(y_i^{(B+1)} | x_i(m, \mathbf{s}^{(B+1)}), s_i^{(B+1)}) > 0, \quad (405)$$

whereas (404) implies for $m' \neq m$

$$\prod_{i=1}^{k'} W(y_i^{(B+1)} | x_i(m', \tilde{\mathbf{s}}), \tilde{s}_i) = 0, \quad (\forall \tilde{\mathbf{s}} \text{ s.t. } (m', \tilde{\mathbf{s}}) \in \mathcal{J}_B). \quad (406)$$

The decoder can thus calculate $\prod_i W(y_i^{(B+1)} | x_i(\tilde{m}, \tilde{\mathbf{s}}), \tilde{s}_i)$ for each $(\tilde{m}, \tilde{\mathbf{s}}) \in \mathcal{J}_B$ and produce the message \tilde{m} for which this product is positive for some $\tilde{\mathbf{s}}$ for which $(\tilde{m}, \tilde{\mathbf{s}}) \in \mathcal{J}_B$.

We next show that, for some choice of k' , there exist codewords $\{\mathbf{x}(m', s')\}$ satisfying (404). To this end we use a random coding argument. Draw the length- k' codewords $\{\mathbf{X}(m', s')\}$ independently, each uniformly over $\mathcal{X}^{k'}$, and let

$$q = \left\lfloor \frac{k'}{l} \right\rfloor, \quad (407a)$$

$$\lambda^* = \min_{s, s' \in \mathcal{S}: \lambda(s) + \lambda(s') > 2\Lambda} \frac{\lambda(s) + \lambda(s')}{2}, \quad (407b)$$

$$\alpha = \frac{\lambda^* - \Lambda}{\lambda^* - \lambda_{\min}}, \quad (407c)$$

$$k'' = \lceil \alpha q l \rceil. \quad (407d)$$

From the cost constraint (396) and the definition of q (407a) it follows that every pair of (not necessarily distinct) state sequences $\mathbf{s}', \mathbf{s}'' \in \mathcal{S}^{k'}$ for which

$$\exists m', m'' \in \mathcal{M} \text{ s.t. } (m', s'), (m'', s'') \in \mathcal{J}_B \quad (408)$$

satisfies

$$\frac{1}{ql} \sum_{i=1}^{ql} \frac{\lambda(s'_i) + \lambda(s''_i)}{2} \leq \Lambda. \quad (409)$$

As we argue next, (409) can hold only if for all such $\mathbf{s}', \mathbf{s}''$ there exist at least k'' distinct epochs $\mathcal{L}(\mathbf{s}', \mathbf{s}'') \subset [1 : ql]$ for which

$$\frac{\lambda(s'_\ell) + \lambda(s''_\ell)}{2} \leq \Lambda, \quad \forall \ell \in \mathcal{L}(\mathbf{s}', \mathbf{s}''). \quad (410)$$

To simplify the typography, we shall refer to $\mathcal{L}(\mathbf{s}', \mathbf{s}'')$ as \mathcal{L} . The claim can then be stated equivalently as

$$\exists \mathcal{L} \subset [1 : ql] \text{ s.t. } (|\mathcal{L}| = k'') \wedge \left(\frac{\lambda(s'_\ell) + \lambda(s''_\ell)}{2} \leq \Lambda, \quad \forall \ell \in \mathcal{L} \right). \quad (411)$$

To prove (411), note that, by the definition of λ^* (407b),

$$\left(\left(\frac{\lambda(s'_i) + \lambda(s''_i)}{2} > \Lambda \right) \implies \left(\frac{\lambda(s'_i) + \lambda(s''_i)}{2} \geq \lambda^* \right) \right), \quad \forall i \in [1 : k'], \quad (412)$$

and, because $\lambda(s) \geq \lambda_{\min}$, $s \in \mathcal{S}$,

$$\frac{\lambda(s'_i) + \lambda(s''_i)}{2} \geq \lambda_{\min}, \quad \forall i \in [1 : k']. \quad (413)$$

The definitions of α (407c) and k'' (407d) combine with (412) and (413) to prove our claim that (411) holds for all $\mathbf{s}', \mathbf{s}'' \in \mathcal{S}^{k'}$ satisfying (408). An immediate consequence of (411) and the assumption (49) is that for all $\mathbf{s}', \mathbf{s}''$ satisfying (408).

$$\begin{aligned} \exists \mathcal{L} \subset [1 : ql] \text{ s.t. } (|\mathcal{L}| = k'') \wedge & \left(\forall \ell \in \mathcal{L} \quad \exists x', x'' \in \mathcal{X} \text{ s.t.} \right. \\ & \left. \left(W(y|x', s'_\ell) W(y|x'', s''_\ell) = 0, \quad \forall y \in \mathcal{Y} \right) \right). \end{aligned} \quad (414)$$

Having established (414), we are now ready to show that—for some choice of k' —the probability that the random codewords $\{\mathbf{X}(m', \mathbf{s}')\}$ satisfy (404) is positive. For every distinct $(m', \mathbf{s}'), (m'', \mathbf{s}'') \in \mathcal{J}_B$

$$\begin{aligned} & \mathbb{P} \left[\forall i \in [1 : k'] \quad \exists y \in \mathcal{Y} \text{ s.t. } W(y|X_i(m', \mathbf{s}'), s'_i) W(y|X_i(m'', \mathbf{s}''), s''_i) > 0 \right] \\ & \leq \left(1 - \frac{1}{|\mathcal{X}|^2} \right)^{k''} \\ & = 2^{-k''(2 \log |\mathcal{X}| - \log(|\mathcal{X}|^2 - 1))}, \end{aligned} \quad (415)$$

$$(416)$$

where we used (414) and that $\mathbf{X}(m', \mathbf{s}')$ and $\mathbf{X}(m'', \mathbf{s}'')$ are independent and uniform over $\mathcal{X}^{k'}$. This, the Union-of-Events bound, and (403) imply that the probability that the randomly drawn length- k' codewords do not satisfy (404) is upper-bounded by

$$\begin{aligned} |\mathcal{J}_B|^2 2^{-k''(2 \log |\mathcal{X}| - \log(|\mathcal{X}|^2 - 1))} \\ \leq 2^{-k''(2 \log |\mathcal{X}| - \log(|\mathcal{X}|^2 - 1)) + 2(k \log |\mathcal{U}| + B \log(1+k) |\mathcal{S}|)}, \end{aligned} \quad (417)$$

which is smaller than one whenever

$$k'' > \frac{k \log |\mathcal{U}| + B \log(1+k) |\mathcal{S}|}{\log |\mathcal{X}| - \frac{1}{2} \log(|\mathcal{X}|^2 - 1)}. \quad (418)$$

Consequently, (407a) and (407d) imply that, if we choose

$$k' = \left(\left\lfloor \frac{k \log |\mathcal{U}| + B \log(1+k) |\mathcal{S}|}{\alpha l (\log |\mathcal{X}| - \frac{1}{2} \log(|\mathcal{X}|^2 - 1))} \right\rfloor + 1 \right) l, \quad (419)$$

then there exist length- k' codewords $\{\mathbf{x}(m', \mathbf{s}')\}$ satisfying (404).

We are now ready to join the dots and conclude that the coding scheme asymptotically achieves any rate smaller than the RHS of (50). More precisely, we will show that, for every rate R smaller than the RHS of (50) and every sufficiently-large blocklength n , our coding scheme can convey nR bits error-free in n channel uses. It follows from (401) and (419) that if the positive integers n , B , k and $\epsilon > 0$ are such that k is a multiple of l ,

$$k \geq \eta_0(|\mathcal{X}|, |\mathcal{S}|, \epsilon), \quad (420a)$$

and

$$\begin{aligned} nR + \left(\left\lfloor \frac{k \log |\mathcal{U}| + B \log(1+k) |\mathcal{S}|}{\alpha l (\log |\mathcal{X}| - \frac{1}{2} \log(|\mathcal{X}|^2 - 1))} \right\rfloor + 1 \right) l \log |\mathcal{S}| \\ \leq Bk \left(\min_{\substack{P_S: \\ \mathbb{E}[\lambda(\tilde{S})] \leq \Lambda}} \max_{P_{U, \mathbf{X}|S}} \min_{\substack{P_{Y|U, \mathbf{X}, S}: \\ P_{Y|U=u, \mathbf{X}, S} \in \mathcal{P}(W), \forall u \in \mathcal{U}}} I(U; Y) - I(S; Y) - \delta(\epsilon, k) \right), \end{aligned} \quad (420b)$$

then the first $B + 1$ blocks of our coding scheme can convey nR bits error-free in

$$Bk + \left(\left\lfloor \frac{k \log |\mathcal{U}| + B \log(1+k) |\mathcal{S}|}{\alpha l (\log |\mathcal{X}| - \frac{1}{2} \log(|\mathcal{X}|^2 - 1))} \right\rfloor + 1 \right) l \quad (421)$$

channel uses. It thus remains to exhibit positive integers B , k , where k is a multiple of l , and some $\epsilon > 0$ such that for every sufficiently-large blocklength n (420) holds and

$$Bk + \left(\left\lfloor \frac{k \log |\mathcal{U}| + B \log(1+k) |\mathcal{S}|}{\alpha l (\log |\mathcal{X}| - \frac{1}{2} \log(|\mathcal{X}|^2 - 1))} \right\rfloor + 1 \right) l \leq n. \quad (422)$$

(When the inequality in (422) is strict, then Block $(B+2)$ deals with all the superfluous epochs: recall that in this block the encoder can transmit arbitrary inputs with the decoder ignoring the corresponding outputs.) As we argue next, when n is sufficiently large we can choose

$$B = \lfloor \sqrt{n} \rfloor - \left(\left\lfloor \frac{\log |\mathcal{U}| + \log(1 + \sqrt{n}) |\mathcal{S}|}{\alpha l (\log |\mathcal{X}| - \frac{1}{2} \log(|\mathcal{X}|^2 - 1))} \right\rfloor + 1 \right) l, \quad (423a)$$

$$k = \left\lfloor \frac{\sqrt{n}}{l} \right\rfloor l, \quad (423b)$$

and we can choose $\epsilon > 0$ for which

$$R + \epsilon < \min_{\substack{P_S: \\ \mathbb{E}[\lambda(S)] \leq \Lambda}} \max_{P_{U,X|S}} \min_{\substack{P_{Y|U,X,S}: \\ P_{Y|U=u,X,S} \in \mathcal{P}(W), \forall u \in \mathcal{U}}} I(U; Y) - I(S; Y). \quad (424)$$

Note that, whenever n is sufficiently large, B and k are positive, k is a multiple of l , and (422) is satisfied. To see that also (420) holds whenever n is sufficiently large, we first observe from (423b) that k tends to infinity as n tends to infinity. This implies that (420a) holds whenever n is sufficiently large, and that $\delta(\epsilon, k)$ (which is defined in (106), where $\gamma_k = \gamma_k(|\mathcal{U}|, |\mathcal{X}|, |\mathcal{S}|, |\mathcal{Y}|)$ converges to zero as k tends to infinity) converges to ϵ as n tends to infinity. We next observe that (423) implies that Bk/n converges to one as n tends to infinity, and that

$$\frac{1}{n} \left(\left\lfloor \frac{k \log |\mathcal{U}| + B \log(1 + k) |\mathcal{S}|}{\alpha l (\log |\mathcal{X}| - \frac{1}{2} \log(|\mathcal{X}|^2 - 1))} \right\rfloor + 1 \right) l \log |\mathcal{S}| \quad (425)$$

converges to zero as n tends to infinity. This, combined with the fact that $\delta(\epsilon, k)$ converges to ϵ as n tends to infinity and with (424), implies that (420b) holds whenever n is sufficiently large. \square

We next prove the converse part of Theorem 2.26.

Converse Part. We first show that (49) is necessary for $C_{f,0}^{(2)}(\Lambda)$ to be positive. To this end suppose that (49) does not hold, i.e., that there exists a pair of states $s, s' \in \mathcal{S}$ satisfying

$$\frac{\lambda(s) + \lambda(s')}{2} \leq \Lambda \quad (426)$$

for which

$$\forall x, x' \in \mathcal{X} \quad \exists y \in \mathcal{Y} \text{ s.t. } W(y|x, s) W(y|x', s') > 0. \quad (427)$$

We will show that in this case it is impossible to transmit a single bit error-free whenever l is even. This will imply that $C_{f,0}^{(2)}(\Lambda, l)$ is zero whenever l is even and consequently that $C_{f,0}^{(2)}(\Lambda)$ is zero, because, by definition,

$$C_{f,0}^{(2)}(\Lambda) = \liminf_{l \rightarrow \infty} C_{f,0}^{(2)}(\Lambda, l). \quad (428)$$

Fix some even l and s, s' as above. The proof is similar to that of the converse of Theorem 2.3. Let the bit take values in the set $\mathcal{M} = \{0, 1\}$, and fix a blocklength n and n encoding mappings

$$f_i: \mathcal{M} \times \mathcal{S}^n \times \mathcal{Y}^{i-1} \rightarrow \mathcal{X}, \quad i \in [1 : n].$$

Denote by $\hat{\mathbf{s}}, \check{\mathbf{s}} \in \mathcal{S}^n$ the state sequences that at odd times are s and s' , respectively, and at even times s' and s , respectively:

$$(\hat{s}_{2i-1}, \check{s}_{2i-1}) = (s, s'), \quad i \in [1 : \lceil n/2 \rceil], \quad (429a)$$

$$(\hat{s}_{2i}, \check{s}_{2i}) = (s', s), \quad i \in [1 : \lfloor n/2 \rfloor]. \quad (429b)$$

Note that, by (426) and because l is even, they meet the cost constraint (48). The line of argument leading to (81) in the converse of Theorem 2.3 implies that there exists an output sequence $\mathbf{y} \in \mathcal{Y}^n$ for which

$$W(y_i | f_i(0, \hat{\mathbf{s}}, y^{i-1}), \hat{s}_i) W(y_i | f_i(1, \check{\mathbf{s}}, y^{i-1}), \check{s}_i) > 0, \quad \forall i \in [1 : n]. \quad (430)$$

This rules out error-free transmission, because if the state sequence is either $\hat{\mathbf{s}}$ or $\check{\mathbf{s}}$, then the decoder, not knowing which, cannot recover the bit.

We next show that—irrespective of whether or not (49) holds— $C_{f,0}^{(2)}(\Lambda)$ is upper-bounded by the RHS of (50). The proof is similar to the converse of Theorem 2.4, but in the current setting we cannot fix some PMF \tilde{P}_S on \mathcal{S} and assume that the state sequence S^n is drawn IID \tilde{P}_S , because this might violate the cost constraint (48). In fact, (48) need not hold even if $\mathbb{E}[\lambda(S)] \leq \Lambda$ under \tilde{P}_S .

Fix any $l \in \mathbb{N}$, and assume that $n = Jl$ for some $J \in \mathbb{N}$. We can make this assumption w.l.g., because

$$\lim_{n \rightarrow \infty} \frac{\lfloor n/l \rfloor l}{n} = 1.$$

To satisfy (48), we fix some l -type \tilde{P}_S on \mathcal{S} w.r.t. which

$$\mathbb{E}[\lambda(S)] \leq \Lambda, \quad (431)$$

and we set \tilde{P}_{S^n} to be the uniform distribution over $(\mathcal{T}_{\tilde{P}_S}^{(l)})^J$. Let the PMF $\tilde{P}_{M,S^n,X^n,Y^n}$ be as in (121) but with \tilde{P}_S^n replaced by \tilde{P}_{S^n} . We can now upper-bound $\frac{1}{n} \log |\mathcal{M}|$ essentially along the line of argument leading to (127) in the converse of Theorem 2.4. The main difference is that under $\tilde{P}_{M,S^n,X^n,Y^n}$ of the current setting S_i and S_{i+1}^n need not be independent and consequently

$$\frac{1}{n} \sum_{i=1}^n I(S_{i+1}^n; S_i)$$

need not be zero. However, it does tend to zero as l tends to infinity, because

$$\begin{aligned} & \frac{1}{n} \sum_{i=1}^n I(S_{i+1}^n; S_i) \\ & \stackrel{(a)}{=} \frac{1}{n} \sum_{i=1}^n [H(S_i) - H(S_i | S_{i+1}^n)] \end{aligned} \quad (432)$$

$$\stackrel{(b)}{=} H(\tilde{P}_S) - \frac{1}{n} H(S^n) \quad (433)$$

$$\stackrel{(c)}{=} H(\tilde{P}_S) - \frac{1}{n} \log |\mathcal{T}_{\tilde{P}_S}^{(l)}|^J \quad (434)$$

$$\stackrel{(d)}{\leq} H(\tilde{P}_S) - \frac{J}{n} (l H(\tilde{P}_S) - \log(1+l) |\mathcal{S}|) \quad (435)$$

$$\stackrel{(e)}{=} \frac{\log(1+l) |\mathcal{S}|}{l} \quad (436)$$

$$\rightarrow 0 \ (l \rightarrow \infty), \quad (437)$$

where (a) holds by the definition of mutual information; (b) follows from the chain rule and the fact that $S_i \sim \tilde{P}_S$ under $\tilde{P}_{M,S^n,X^n,Y^n}$; (c) holds because S^n is uniform over $(\mathcal{T}_{\tilde{P}_S}^{(l)})^J$ under $\tilde{P}_{M,S^n,X^n,Y^n}$; (d) follows from the inequality $|\mathcal{T}_{\tilde{P}_S}^{(l)}| \geq (1+l)^{-|\mathcal{S}|} 2^{lH(\tilde{P}_S)}$, where the entropy is computed w.r.t. \tilde{P}_S [8, Lemma 2.3]; and (e) holds because $n = Jl$.

Having established (437), we are now ready to conclude the proof. The arguments leading to (139) in the converse of Theorem 2.4 and (436) imply that

$$\begin{aligned} \frac{1}{n} \log |\mathcal{M}| & \leq \frac{\log(1+l) |\mathcal{S}|}{l} \\ & + \max_{\tilde{P}_{U,X|S}} \min_{\substack{\tilde{P}_{Y|U,X,S}: \\ \tilde{P}_{Y|U=u,X,S} \in \mathcal{P}(W), \forall u \in \mathcal{U}}} I(U; Y) - I(U; S), \end{aligned} \quad (438)$$

where U is an auxiliary chance variable taking values in a finite set \mathcal{U} , and the mutual informations are computed w.r.t. the joint PMF $\tilde{P}_S \times \tilde{P}_{U,X|S} \times \tilde{P}_{Y|U,X,S}$. Moreover, it is enough to consider the second term on the RHS of (438), because the first converges to zero as l tends to infinity (437) and, by definition,

$$C_{f,0}^{(2)}(\Lambda) = \liminf_{l \rightarrow \infty} C_{f,0}^{(2)}(\Lambda, l).$$

To conclude that $C_{f,0}^{(2)}(\Lambda)$ is upper-bounded by the RHS of (50), we would have liked to choose some PMF \tilde{P}_S that—among all PMFs on \mathcal{S} w.r.t. which (431) holds—yields the tightest bound, i.e., minimizes

$$\max_{\tilde{P}_{U,X|S}} \min_{\substack{\tilde{P}_{Y|U,X,S}: \\ \tilde{P}_{Y|U=u,X,S} \in \mathcal{P}(W), \forall u \in \mathcal{U}}} I(U; Y) - I(U; S). \quad (439)$$

But this is not possible, because \tilde{P}_S must be an l -type. We can, however, choose \tilde{P}_S as one that—among all l -types on \mathcal{S} w.r.t. which (431) holds—minimizes (439). For this choice (438) implies that

$$C_{f,0}^{(2)}(\Lambda) \leq \liminf_{l \rightarrow \infty} \min_{\substack{\tilde{P}_S \in \Gamma^{(l)}: \\ \mathbb{E}[\lambda(s)] \leq \Lambda}} \max_{\tilde{P}_{U,X|S}} \min_{\substack{\tilde{P}_{Y|U,X,S}: \\ \tilde{P}_{Y|U=u,X,S} \in \mathcal{P}(W), \forall u \in \mathcal{U}}} I(U; Y) - I(U; S), \quad (440)$$

where $\Gamma^{(l)}$ denotes the set of l -types on \mathcal{S} . To conclude, note that the RHS of (440) is equal to that of (50): every PMF \tilde{P}_S on \mathcal{S} w.r.t. which (431) holds can be approximated in the total variation distance by an l -type on \mathcal{S} w.r.t. which (431) holds when l is sufficiently large; and (conditional) entropy is continuous in this distance [8, Lemma 2.7]. \square

References

- [1] C. E. Shannon, “The zero error capacity of a noisy channel,” *IRE Trans. Inf. Theory*, vol. 2, no. 3, pp. 8–19, Sep. 1956.
- [2] S. I. Gel’fand and M. S. Pinsker, “Coding for channel with random parameters,” *Problems of Control Theory*, vol. 9, no. 1, pp. 19–31, 1980.
- [3] N. Merhav and T. Weissman, “Coding for the feedback Gel’fand-Pinsker channel and the feedforward Wyner-Ziv source,” *Proc. of IEEE Int. Symp. on Inf. Theory (ISIT)*, pp. 1506–1510, Sep. 2005.
- [4] G. Dueck, “The zero error feedback capacity region of a certain class of multiple-access channels,” *Problems of Control and Inf. Theory*, vol. 14, no. 2, pp. 89–103, 1985.
- [5] R. Ahlswede, “Channels with arbitrarily varying channel probability functions in the presence of noiseless feedback,” *Zeitschrift f. Wahrscheinlichkeitstheorie und verw. Gebiete*, vol. 25, no. 3, pp. 239–252, Sep. 1973.

- [6] J. M. Ooi and G. W. Wornell, “Fast iterative coding techniques for feedback channels,” *IEEE Trans. Inf. Theory*, vol. 44, no. 7, pp. 2960–2976, Nov. 1998.
- [7] R. Ahlswede, “A note on the existence of the weak capacity for channels with arbitrarily varying channel probability functions and its relation to Shannon’s zero error capacity,” *Ann. of Math. Stat.*, vol. 41, no. 3, pp. 1027–1033, Jun. 1970.
- [8] I. Csiszár and J. Körner, *Information Theory*, 2nd. ed., Cambridge University Press, 2011.
- [9] R. Ahlswede, “Arbitrarily varying channels with states sequence known to the sender,” *IEEE Trans. Inf. Theory*, vol. 32, no. 5, pp. 621–629, Sep. 1986.
- [10] C. E. Shannon, “A mathematical theory of communication,” *The Bell System Tech. J.*, vol. 27, pp. 379–423 and 626–656, July and Oct. 1948.
- [11] C. E. Shannon, “Channels with side-information at the transmitter,” *IBM J. of Research and Develop.*, vol. 2, pp. 289–293, Oct. 1958.
- [12] A. El Gamal and Y.-H. Kim, *Network Information Theory*, Cambridge University Press, 2011.
- [13] Y.-H. Kim and A. Sutivong and T. M. Cover, “State Amplification,” *IEEE Trans. Inf. Theory*, vol. 54, no. 5, pp. 1850–1859, May 2008.
- [14] C. Choudhuri and Y.-H. Kim and U. Mitra, “Causal State Communication,” *IEEE Trans. Inf. Theory*, vol. 59, no. 6, pp. 3709–3719, Jun. 2013.
- [15] S. Bross and A. Lapidoth, “Conveying Data and State with Feedback,” *to appear in Proc. of IEEE Int. Symp. on Inf. Theory (ISIT)*, Jul 2016.
- [16] F. M. J. Willems and E. C. Van Der Meulen, “The discrete memoryless multiple-access channel with cribbing encoders,” *IEEE Trans. Inf. Theory*, vol. 31, no. 3, pp. 313–327, May 1985.